

# When Does Geometry Emerge from Memorization in Transformers?

Anonymous authors  
Paper under double-blind review

## Abstract

Transformer models often show structured internal representations on relational tasks, which are often interpreted as geometric organization. Prior work documents such structure via visualization or performance-based analyses, but does not isolate whether perfect memorization alone yields geometric representations. Here, we conduct a controlled study using synthetic relational worlds defined by canonical graph topologies, explicitly training Transformers to perfectly memorize relational structure without imposing constraints that favor or discourage geometry, and examining whether geometric representations arise as a consequence.

Across chains, cycles, regular graphs, and star graphs, models achieve perfect memorization accuracy while internal embeddings do not systematically preserve either global distances or local neighborhoods, indicating reliance on non-geometric, index-based representations. By probing embeddings against shortest-path distance using rank consistency and neighborhood preservation metrics, we show that memorization alone places no requirement on metric organization. Recoverable geometric structure emerges only when the task objective, together with the relational topology, sufficiently constrains node interchangeability, reducing the space of symmetry-equivalent memorization solutions.

Our results show that perfect memorization does not imply emergent geometric structure, and characterize the conditions under which structure arises in learned embeddings.

## 1 Introduction

Transformer models are frequently observed to exhibit structured internal representations, such as linear arrangements or clustered layouts, which are often interpreted as evidence that the model has learned a geometric representation of underlying relations (Mikolov et al., 2013b; Hewitt & Manning, 2019). Such observations have motivated the view that geometric organization may emerge naturally during training, and are commonly taken as evidence of relational structure encoded in neural representations (Petroni et al., 2019).

At the same time, achieving perfect accuracy on a relational task does not require preserving distances, neighborhoods, or any notion of metric structure. Prior work has shown that neural models can achieve perfect memorization through arbitrary internal coding schemes, such as index-based or identity-based representations, without analyzing whether such solutions induce any coherent relational or geometric structure (Zhang et al., 2017; Arpit et al., 2017). Transformer architectures in particular admit such associative or lookup-style solutions (Geva et al., 2021), making high task performance alone insufficient to diagnose whether learned representations are genuinely geometric.

This ambiguity complicates the interpretation of probe-based geometric analyses of learned representations: apparent structure may arise even when representations are not uniquely identifiable or causally required for task performance (Hewitt et al., 2021; Elazar et al., 2021). In particular, low-dimensional visualizations from techniques such as UMAP can appear highly structured despite substantial distortions in the underlying metric, and therefore cannot, by themselves, establish meaningful geometric organization (McInnes et al.,

2018). Accordingly, we adopt an operational definition of geometric organization based on consistency between distances in representation space and shortest-path distances in the underlying relational graph.

We study the relationship between memorization, internal geometric organization, and representational stability in a controlled setting. Transformer models are trained from scratch on synthetic relational graph worlds under objectives that permit perfect memorization without requiring generalization, allowing us to isolate when geometric structure is forced by the task rather than incidental. We show that memorization alone does not reliably induce geometric organization, and that stable geometry emerges only when the task objective and relational topology sufficiently constrain node interchangeability, reducing the space of symmetry-equivalent memorization solutions.

## 2 Methodology

### 2.1 Objective and Scope

We investigate whether perfect memorization of relational structure necessarily induces a geometric internal representation. To isolate this question, we study Transformer models trained under a pure memorization regime, where the learning objective can be satisfied without compositional reasoning, generalization, or abstraction. Our analysis focuses on the geometry of internal representations rather than task performance alone.

All experiments are conducted in closed synthetic environments, allowing us to disentangle the effects of topology, model capacity, and embedding dimensionality without confounds from linguistic priors or pre-training.

### 2.2 Why Graphs Rather Than Lexical Data

We analyze synthetic graph-structured data rather than lexical or natural language inputs to achieve precise experimental control. Lexical data introduce confounding factors such as semantics, frequency effects, distributional similarity, and pretraining priors, which obscure the relationship between memorization and internal representation.

Graphs provide an explicit relational structure with a well-defined notion of distance, symmetry, and topology. This allows direct comparison between internal embedding geometry and ground-truth relational distances, enabling us to attribute observed representational effects to structural properties of the task rather than linguistic artifacts.

### 2.3 Relational Graph Worlds

Each experiment is defined by a graph  $G = (V, E)$ , where nodes represent discrete entities and edges represent valid successor relations. Given a node  $u \in V$ , the model is trained to predict one of its successors  $v \in \mathcal{N}(u)$ , where  $\mathcal{N}(u)$  denotes the set of valid successors of  $u$ . For chain, cycle, star, and 3-regular graphs, successor relations are directed, ensuring that successor prediction is unambiguous and reducing additional symmetry-based degeneracies; for Watts–Strogatz graphs, edges are undirected and represented using bidirectional successor relations.

Throughout this work, we use *symmetry* to refer specifically to automorphisms of the underlying graph. Graph automorphisms determine when nodes are structurally interchangeable under relabeling. However, node interchangeability under the training objective may also arise from task-induced equivalence, even in graphs with limited formal symmetry. In both cases, large equivalence classes reduce the identifiability of metric structure, permitting multiple functionally correct but geometrically distinct representations. We consider several canonical graph families, each selected to probe a distinct structural property relevant to the emergence of geometric representations.

- 80 • **Chain.** A directed path with a unique successor per node. Although the graph has a clear linear  
81 ordering, the successor-prediction objective does not require the model to represent this ordering ex-  
82 plicitly, allowing multiple functionally equivalent encodings that need not preserve metric structure.
- 83 • **Cycle.** A closed loop with uniform degree and rotational symmetry, making all nodes structurally  
84 interchangeable and globally ambiguous.
- 85 • **Star.** A hub-and-spoke topology with one central node and many structurally identical leaves,  
86 yielding shallow and degenerate graph distances.
- 87 • **3-Regular Graph.** A graph where each node has degree three. The topology is locally uniform  
88 but does not impose a unique global ordering, reducing symmetry without introducing privileged  
89 nodes.
- 90 • **Watts–Strogatz Small-World Graph.** A graph combining local regularity with randomly  
91 rewired long-range connections, introducing multiple length scales and partial symmetry breaking.

92 Unless otherwise stated, graphs have  $|V| = 60$  nodes; we also test  $|V| \in \{15, 150\}$  to assess finite-size and  
93 scaling effects.

## 94 2.4 Model Architecture and Training Regime

95 We use a decoder-only Transformer trained from scratch, without pretraining or auxiliary objectives. Each  
96 node  $u$  is mapped to a unique discrete token. Given a single node token the model predicts a distribution  
97 over node tokens for valid successors. Although we instantiate a decoder-only Transformer, inputs consist of  
98 a single token, so the training objective does not engage sequence-level dependencies or multi-step autore-  
99 gressive structure. To verify that our findings are not artifacts of the decoder-only architecture, we replicate  
100 our core experiments using an encoder-decoder framework, which yields identical conclusions (Appendix C).

101 Models are trained with cross-entropy (CE) loss for successor prediction, which permits assigning high  
102 probability to a single valid successor without enforcing uniformity across neighbors. We also ran a small-  
103 scale comparison using a  $1/k$  uniform target over neighbors (where  $k = |N(u)|$ ) and observed no qualitative  
104 change in the resulting internal geometry. We report CE-based results throughout, as it is the standard and  
105 least-constraining objective for this setting.

106 Training uses a fixed number of epochs sufficient to reach stable convergence and near-perfect set accuracy  
107 ( $\geq 99.9\%$ ) across graph families. In deterministic topologies such as chains and cycles, the cross-entropy loss  
108 approaches zero. In topologies with multiple valid successors, loss stabilizes at a non-zero value determined  
109 by the successor distribution, approaching but not necessarily reaching the uniform-entropy bound, as the  
110 training objective does not enforce uniformity across valid successors.

## 111 2.5 Embedding Dimension and Graph Size Sweeps

112 To disentangle geometric effects from representational capacity, we vary both embedding dimensionality and  
113 graph size.

114 Embedding dimension  $d$  is varied (e.g., 32, 64, 128, 256), testing whether geometric organization emerges from  
115 increased capacity. Graph size  $|V|$  is varied across 15, 60, and 150 nodes, probing whether larger relational  
116 worlds constrain representations more strongly or permit non-geometric memorization strategies to persist.

117 All other architectural and optimization settings are held fixed across these sweeps.

## 118 2.6 Verification of Memorization

119 To verify that models successfully memorize the relational structure, we compute *set accuracy*. For each  
120 node  $u$ , a prediction is counted as correct if the model’s argmax prediction lies in  $\mathcal{N}(u)$ . Across all reported  
121 settings, models achieve set accuracy above 99.9%, ensuring that differences in internal representations are  
122 not attributable to task failure.

## 123 2.7 Geometric Probing of Internal Representations

124 To assess whether memorization corresponds to geometric organization, we analyze frozen hidden representa-  
 125 tions at every layer. Let  $d_G(u, v)$  denote the shortest-path distance between nodes  $u$  and  $v$  in the graph, and  
 126 let  $d_E(u, v) = \|e_u - e_v\|_2$  denote the Euclidean distance between their embedding vectors. Graph distances  
 127 are computed using shortest paths; node pairs that are unreachable under the graph are excluded from all  
 128 distance-based probes.

### Spearman Rank Correlation.

$$\rho = \text{Spearman}(d_G(u, v), d_E(u, v)) \quad (1)$$

129 This measures whether embedding distances preserve the ordering of graph distances. High values indicate  
 130 that closer nodes in the graph tend to be closer in the embedding space, even if absolute scaling is distorted.

### Coefficient of Determination ( $R^2$ ).

$$R^2 = 1 - \frac{\sum (d_E - \hat{d}_E)^2}{\sum (d_E - \bar{d}_E)^2}, \quad \hat{d}_E = a d_G + b \quad (2)$$

131 Here,  $\hat{d}_E$  denotes the *least-squares linear fit* of embedding distances  $d_E(u, v)$  to graph shortest-path distances  
 132  $d_G(u, v)$  over all node pairs  $(u, v)$ , and  $\bar{d}_E$  denotes the mean embedding distance. This measures how well  
 133 graph distances linearly predict embedding distances; high values indicate proportional alignment consistent  
 134 with a global metric embedding.

### Metric Distortion.

$$\text{Distortion} = \frac{\max_{u \neq v} \frac{d_E(u, v)}{d_G(u, v)}}{\min_{u \neq v} \frac{d_E(u, v)}{d_G(u, v)}} \quad (3)$$

135 This quantifies the uniformity of distance scaling in a worst-case sense. Low distortion indicates the absence  
 136 of extreme pairwise scaling deviations, while high distortion reflects uneven stretching or compression.

### Expansion Variance.

$$\text{Var}\left(\frac{d_E(u, v)}{d_G(u, v)}\right) \quad (4)$$

137 Expansion variance measures how uniformly graph distances are scaled in embedding space. If a coherent  
 138 global metric is recovered, distances at different graph scales are expanded by similar factors, resulting in  
 139 relatively low variance. Topologies that exhibit asymmetry, boundary effects, or anisotropy can exhibit non-  
 140 uniform scaling even when geometry is well formed. In contrast, graphs with degenerate distance structure  
 141 can produce artificially low variance because many node pairs share identical graph distances.

### Local Isometry (1-Hop Stretch).

$$\frac{1}{|V|} \sum_{u \in V} \frac{1}{|\mathcal{N}(u)|} \sum_{v \in \mathcal{N}(u)} d_E(u, v) \quad (5)$$

142 Local isometry measures how consistently immediate graph neighborhoods are represented in embedding  
 143 space. When geometry is well formed, neighboring nodes are embedded at comparable distances, indicating  
 144 stable local structure. Expected values depend on topology, with shallow neighborhoods yielding smaller  
 145 values and richer local connectivity producing larger values. This metric is therefore interpreted relative to  
 146 graph structure rather than as an absolute criterion.

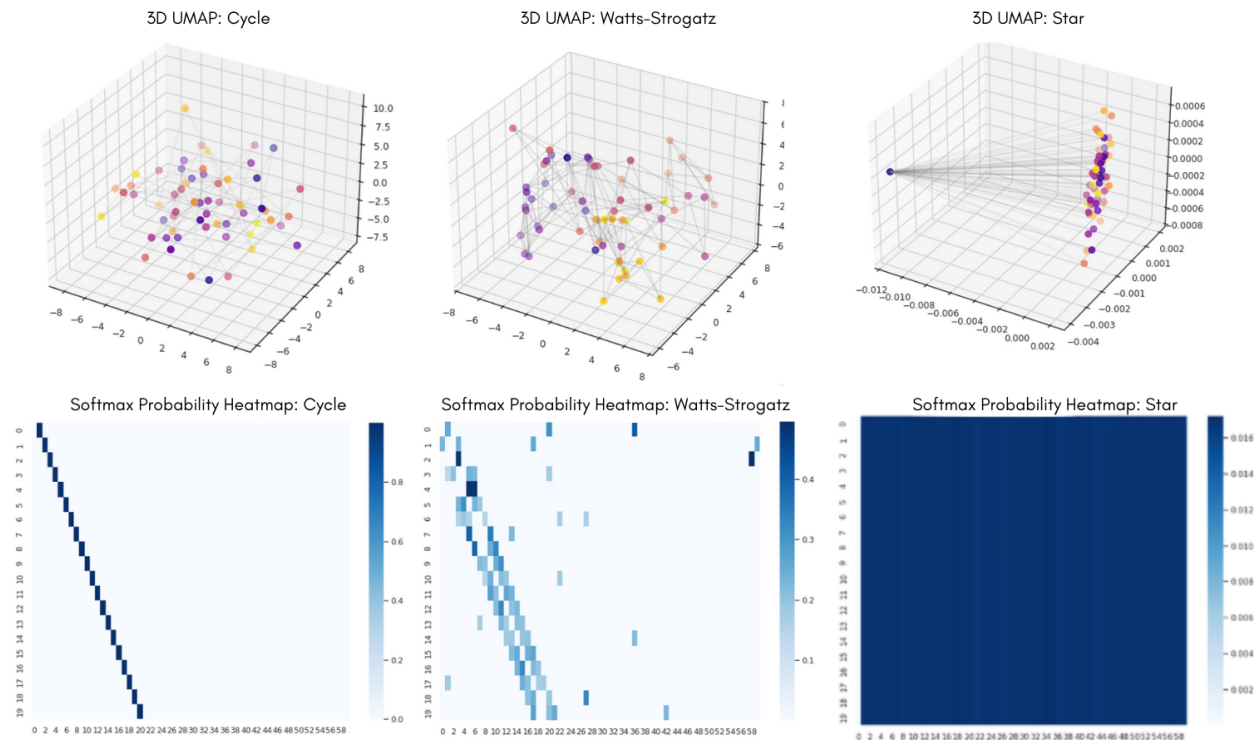


Figure 1: **Geometry vs. memorization across graph topologies.** Top: 3D projections of learned node embeddings. Bottom: successor probability heatmaps. While the model perfectly memorizes local transitions in all cases (bottom), partially recoverable geometric organization emerges only for topologies with sufficient asymmetry (e.g., Watts–Strogatz), and collapses under hub–spoke structure (Star).

## 147 2.8 Interpretation of Geometric Outcomes

148 A model is considered to exhibit *recoverable geometric organization* if embedding distances reflect graph  
 149 distances under our probes, as indicated by high Spearman ( $\rho$ ), high  $R^2$ , controlled distortion, and *topology-*  
 150 *consistent* expansion variance. Local isometry is used as a complementary measure of neighborhood coherence  
 151 and interpreted comparatively across graph families. Conversely, a model is considered non-geometric if it  
 152 achieves perfect memorization while failing to preserve globally or locally consistent metric structure.

## 153 3 Results

154 We evaluate whether geometric structure necessarily emerges from memorization by analyzing internal repre-  
 155 sentations of decoder-only Transformers trained on synthetic relational graphs. Results are reported for  
 156 graphs with  $n = 60$  nodes and embedding dimension  $d = 128$ , a setting that is large enough to admit  
 157 non-trivial relational structure while remaining small enough to allow detailed geometric analysis; additional  
 158 experiments across node counts and embedding dimensions are provided in the Appendix.

### 159 3.1 Perfect Memorization Does Not Imply Geometry

160 Despite identical training objectives and perfect memorization accuracy across all graph families, internal  
 161 representations differ in their geometric properties. Final-layer geometric probe results reveal that memo-  
 162 rization alone is insufficient to guarantee the emergence of a coherent metric structure.

Graph	Spearman $\rho$	$R^2$	Distort.	Exp. Var	Loc. Iso
Cycle	$0.030 \pm 0.027$	$0.001 \pm 0.001$	$70.90 \pm 5.23$	$16.20 \pm 1.09$	$20.35 \pm 0.72$
Chain	$0.028 \pm 0.031$	$0.001 \pm 0.002$	$74.80 \pm 7.04$	$15.71 \pm 1.72$	$19.98 \pm 1.15$
3-Regular	$0.359 \pm 0.094$	$0.103 \pm 0.043$	$5.11 \pm 0.65$	$18.17 \pm 1.13$	$14.23 \pm 0.61$
Watts–Strogatz	$0.673 \pm 0.027$	$0.502 \pm 0.030$	$8.63 \pm 0.59$	$5.37 \pm 0.21$	$12.11 \pm 0.18$
Star	N/A*	N/A*	$6.18 \pm 1.57$	$0.01 \pm 0.00$	$0.18 \pm 0.07$

Table 1: Aggregate final-layer geometric probe results and functional stability across fifteen random seeds ( $n = 60, d = 128$ ). Results report Mean  $\pm$  SD. \*Spearman  $\rho$  and  $R^2$  are mathematically undefined for Star graphs due to the lack of variance in ground-truth shortest-path, rendering correlation and linear regression denominators zero.

### 3.2 Graphs with High Node Interchangeability Collapse to Non-Geometric Memorization

Quantitative probe results for all graph families are reported in Table 1.

**Chain and Cycle.** For both the chain and cycle graphs, Spearman rank correlation between embedding distances and graph distances remains near zero, and distortion values are extremely high. Elevated expansion variance in the absence of rank alignment indicates non-uniform scaling unrelated to graph distance, consistent with projection artifacts rather than meaningful geometry.

As shown in Figure 1 (column 1), 3D projections of cycle form scattered point clouds without global ordering. Nevertheless, the corresponding softmax probability distributions are sharply peaked, confirming that successor prediction is memorized precisely. These results indicate that relational worlds with high node interchangeability permit index-based or lookup-style representations that do not require geometry.

### 3.3 Partial symmetry-breaking: fragile / unstable geometry

**3-Regular Graph.** The 3-regular graph introduces limited asymmetry while preserving uniform degree. Embeddings show moderate rank correlation, indicating partial alignment between embedding and graph distances. While distortion shows no extreme failures, expansion variance and local isometry are inconsistent with the global structure supported by this topology, reflecting non-uniform scaling and unstable neighborhood structure. As a result, global ordering remains weak. These results indicate that limited asymmetry permits weak but incomplete geometric alignment, insufficient to support a stable global metric embedding.

### 3.4 Stronger symmetry-breaking: recoverable geometry

**Watts–Strogatz Graph.** The Watts–Strogatz small-world graph reduces global symmetry while preserving local connectivity. Embeddings achieve higher rank correlation and substantially higher  $R^2$ , indicating stronger recovery of global distance ordering, with moderate distortion and relatively low, topology-consistent expansion variance. Embeddings in Figure 1 (column 2) form a coherent low-dimensional structure in which graph distances are partially recoverable, and successor probabilities assign mass to multiple valid successors due to the presence of multiple admissible paths. Together, these results show that stronger symmetry breaking constrains memorization sufficiently for geometry to become recoverable under our probes.

### 3.5 Star Graph as a Structural Control

The star graph is a control case in which geometric structure is degenerate and imposed by the graph topology. Rank-based metrics are undefined due to the collapse of graph distances, while distortion, local isometry, and expansion variance take trivially non-informative values due to distance degeneracy.

As shown in Figure 1 (column 3), the hub–spoke topology enforces a simple geometric separation between the hub and spokes. This confirms that the probes distinguish geometry imposed by topology from geometric organization that arises under relational constraints.

### 195 3.6 Summary

196 Across all experiments, we observe a consistent pattern: memorization is compatible with multiple represen-  
 197 tational strategies, only some of which are geometric. Relational worlds with large symmetry-equivalence  
 198 classes permit non-geometric solutions, partial symmetry-breaking yields fragile geometry, and sufficiently  
 199 constrained topologies enable stable recoverable geometric organization. Together, these results show that  
 200 geometric structure is not an automatic consequence of memorization, but instead depends critically on the  
 201 topological and symmetry constraints of the underlying relational environment.

## 202 4 Discussion

### 203 4.1 Memorization Permits Multiple Representational Realizations

204 Our experiments show that perfect memorization does not uniquely determine the geometric organization of  
 205 internal representations. Across all graph families, models achieve near-perfect set accuracy and approach  
 206 their task-dependent entropy lower bound, yet exhibit markedly different geometric properties when probed.  
 207 This reflects that the training objective admits multiple internal realizations that are functionally equivalent  
 208 for memorization.

209 Whether geometric structure becomes recoverable depends on how relational constraints are jointly imposed  
 210 by topology and task. Topology-induced constraints arise from structural asymmetries in the graph itself,  
 211 while task-induced constraints arise from the successor-prediction objective, which may render distinct nodes  
 212 functionally interchangeable. Graphs such as Chains and Cycles possess well-defined topological structure,  
 213 but the task objective does not expose this structure, permitting non-geometric memorization solutions. The  
 214 3-regular graph introduces partial asymmetry, weakly restricting interchangeability and inducing fragile, lo-  
 215 cally coherent geometry that does not stabilize into a global metric. Conversely, the Star graph imposes  
 216 strong topological asymmetry, but the task remains degenerate due to collapsed distance structure, yield-  
 217 ing only trivial, topology-imposed geometry. In contrast, the Watts–Strogatz graph introduces persistent  
 218 structural heterogeneity that, when coupled with the successor-prediction objective, restricts the space of ad-  
 219 missible memorization solutions enough for graph distances to become partially recoverable from embedding  
 220 distances under our probes.

### 221 4.2 Symmetry Suppresses Metric Identifiability

222 Here we formalize symmetry in terms of graph automorphisms, which determine when nodes are interchange-  
 223 able under the training objective. Graph families such as cycles admit large automorphism groups, while  
 224 others such as chains exhibit task-induced interchangeability under the successor objective despite limited  
 225 formal symmetry. From a representational perspective, this implies that multiple embeddings related by  
 226 permutation or rotation are equally valid solutions to the training objective. As a consequence, there is no  
 227 unique or preferred metric ordering that must be recovered in embedding space.

228 This symmetry manifests empirically as near-zero rank correlation between embedding distances and graph  
 229 distances, along with persistently high distortion. Because successor-preserving relabelings, whether arising  
 230 from graph automorphisms or from task-induced equivalence, do not affect prediction accuracy, embeddings  
 231 need not align with graph-theoretic distances.

### 232 4.3 Symmetry-Breaking Constrains Representational Degrees of Freedom

233 Graph topologies that introduce asymmetry or irregular connectivity reduce the space of functionally equiv-  
 234 alent representations. When nodes differ in their relational roles, automorphism-induced equivalence classes  
 235 collapse, and representational solutions must respect these distinctions to satisfy the training objective.

236 In such settings, embedding distances become more strongly coupled to graph distances, resulting in im-  
 237 proved rank consistency and reduced distortion. However, this coupling is conditional rather than automatic.  
 238 Increasing embedding dimension expands the space of admissible representations and may suppress geomet-  
 239 ric alignment unless relational asymmetry sufficiently constrains non-metric solutions. This explains why

240 geometry does not emerge monotonically with capacity, and why scale alone does not guarantee metric  
241 structure.

#### 242 **4.4 Capacity, Scale, and Robustness Across $n$ and $d$**

243 Varying the number of nodes and embedding dimension confirms that the observed behaviors are robust  
244 to scale. In symmetric graph families, increasing  $n$  or  $d$  preserves near-perfect memorization while leaving  
245 geometric probes unchanged, indicating that scale does not resolve metric non-identifiability. In contrast, for  
246 asymmetric graph families, increased capacity stabilizes geometric organization when topological constraints  
247 limit alternative realizations.

248 These findings indicate that geometric structure arises from the interaction between model capacity and  
249 relational constraints, rather than from scale or optimization dynamics alone.

#### 250 **4.5 Boundary Cases and Implications**

251 The star graph provides a boundary case in which geometric structure is imposed directly by topology  
252 rather than resolved by the task. Owing to its hub–spoke organization, geometry is trivial and unavoidable  
253 and rank-based probes are degenerate by construction. This illustrates that geometric organization may be  
254 imposed, recoverable, or absent depending on the relational structure.

255 More broadly, these findings indicate that geometric analyses of learned representations should be interpreted  
256 relative to task identifiability, rather than treated as intrinsic properties of the architecture. When the  
257 training objective admits multiple functionally equivalent solutions, the presence or absence of geometric  
258 structure alone is insufficient to diagnose how relational information is encoded.

## 259 **5 Related Work**

### 260 **5.1 Geometric Structure in Transformer Representations**

261 A growing literature reports that Transformer models often exhibit structured or geometric organization  
262 in their hidden representations. Early evidence of linear and hierarchical structure in neural embeddings  
263 emerged from studies of word representations (Mikolov et al., 2013a; Pennington et al., 2014), followed by  
264 analyses showing that contextual Transformers encode syntactic and semantic relations in linearly recoverable  
265 form (Hewitt & Manning, 2019; Jawahar et al., 2019). More recent work has studied the geometric structure  
266 of transformer representations across layers and tasks (Valeriani et al., 2023).

### 267 **5.2 Associative Memory, Lookup Tables, and Memorization**

268 A long-standing view models neural memorization as associative or lookup-based storage of co-occurrences,  
269 where embeddings act as arbitrary identifiers and relations are stored in weight matrices (Hopfield, 1982).  
270 This abstraction has been widely adopted in analyses of Transformer memory (Geva et al., 2021), and is  
271 often sufficient to explain behavior on disjoint or atomic facts.

272 Several works argue that, in the absence of structural constraints, Transformers can implement near-perfect  
273 lookup-table solutions (Arpit et al., 2017). Recent empirical studies suggest that increased model capacity  
274 primarily improves memorization fidelity rather than inducing structured representations. Our findings  
275 align with this view in symmetric relational worlds, while clarifying the conditions under which associative  
276 solutions coexist with or give way to geometric ones.

### 277 **5.3 Probing, Identifiability, and Representation Analysis**

278 Probing methods are commonly used to assess what information is encoded in neural representations (Adi  
279 et al., 2017; Conneau et al., 2018). However, multiple works caution that probe success does not imply that a  
280 representation intrinsically encodes a property, as probes may exploit incidental correlations or task structure  
281 (Hewitt et al., 2021; Ravichander et al., 2021). These observations motivate metric, rank, and local-structure

282 probes beyond linear separability (Ethayarajh, 2019), alongside information-theoretic alternatives (Pimentel  
 283 et al., 2020). Our results highlight a limitation of geometric probing analyses commonly used in natural  
 284 language processing. When the training objective admits multiple functionally equivalent representations,  
 285 observed geometric structure may be incidental rather than identifiable. Probing results should therefore be  
 286 interpreted relative to the constraints imposed by the task, rather than assumed to reflect learned relational  
 287 understanding.

288 Our approach builds on this perspective by evaluating geometric organization using complementary distance-  
 289 based probes and by analyzing robustness across scales. We connect geometric probing to task identifiability:  
 290 when multiple representational realizations are functionally equivalent, probes may reveal the absence of a  
 291 uniquely recoverable geometry.

## 292 5.4 Relational and Graph-Based Learning

293 Relational reasoning has traditionally been studied using graph neural networks, which impose explicit  
 294 inductive biases for locality and message passing (Kipf & Welling, 2017; Hamilton et al., 2017). Transformers,  
 295 despite lacking such biases, have shown strong empirical performance on graph and relational tasks (Dwivedi  
 296 et al., 2023). Several works explore hybrid or geometry-aware architectures that explicitly encourage metric  
 297 structure (Chen et al., 2022).

298 Our work differs from prior work by holding architecture fixed and varying relational topology. This allows  
 299 us to isolate how symmetry and asymmetry in the task constrain representational geometry, independent of  
 300 architectural enforcement.

## 301 6 Conclusion

302 We investigated whether geometric organization is a necessary consequence of memorization in Transformer  
 303 models by analyzing internal representations learned on synthetic relational graph worlds. Across a range  
 304 of graph topologies, node counts, and embedding dimensions, we show that perfect memorization is com-  
 305 patible with multiple representational realizations, only some of which exhibit coherent geometric structure.  
 306 Our results identify topology induced symmetry breaking in the task as a central factor governing when  
 307 geometry is recoverable from the training objective. Relational worlds with large equivalence classes arising  
 308 from graph symmetry or from task-induced interchangeability admit non-geometric memorization strategies  
 309 even at scale. Importantly, increasing model capacity alone does not guarantee geometric emergence and  
 310 may instead expand the space of functionally equivalent, non-metric solutions. These findings clarify the  
 311 distinction between memorization and representation, and caution against interpreting geometric structure  
 312 as an intrinsic or automatic property of Transformer architectures. More broadly, our work highlights the  
 313 role of task identifiability in representation analysis and provides a principled framework for evaluating when  
 314 geometric organization reflects genuine relational structure rather than functionally equivalent alternatives.

## 315 7 Limitations

316 Our study focuses on controlled synthetic relational environments designed to isolate the representational  
 317 consequences of memorization. While this setting enables precise analysis of identifiability and geometric  
 318 structure, it does not capture the full complexity of natural language or multimodal data, where semantic  
 319 grounding, redundancy, and compositional generalization may introduce additional constraints on represen-  
 320 tations.

321 We analyze a fixed Transformer architecture trained under a successor-prediction objective. Different ar-  
 322 chitectural biases, training objectives, or regularization schemes may alter the space of admissible represen-  
 323 tational solutions. In addition, although we employ multiple complementary geometric probes, all probing  
 324 analyses depend on the choice of metric and evaluation procedure. Our conclusions therefore concern the ex-  
 325 istence and stability of *recoverable* geometric structure under these probes, rather than asserting the absence  
 326 of any latent organization.

## References

- 327
- 328 Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. Fine-grained analysis of  
329 sentence embeddings using auxiliary prediction tasks. In *International Conference on Learning Representations*, 2017.  
330
- 331 Devansh Arpit, Stanislaw Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S. Kan-  
332 wal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, and Simon Lacoste-Julien. A closer  
333 look at memorization in deep networks. In *International Conference on Machine Learning*, 2017.
- 334 Dexiong Chen, Leslie O’Bray, and Karsten Borgwardt. Structure-aware transformers for graph representation  
335 learning. In *International Conference on Machine Learning*, 2022.
- 336 Alexis Conneau, Germán Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. What you can  
337 cram into a single vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the*  
338 *56th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2018.
- 339 Vijay Prakash Dwivedi, Chaitanya K. Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier  
340 Bresson. Benchmarking graph neural networks. *Journal of Machine Learning Research*, 2023.
- 341 Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. Amnesic probing: Behavioral explanation  
342 with amnesic counterfactuals. *Transactions of the Association for Computational Linguistics*, 2021.
- 343 Kawin Ethayarajh. How contextual are contextualized word representations? comparing the geometry of  
344 bert, elmo, and gpt-2 embeddings. In *Proceedings of EMNLP-IJCNLP*, 2019.
- 345 Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value  
346 memories. In *Empirical Methods in Natural Language Processing*, 2021.
- 347 William Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In  
348 *Advances in Neural Information Processing Systems*, 2017.
- 349 John Hewitt and Christopher Manning. A structural probe for finding syntax in word representations. In  
350 *Proceedings of NAACL*, 2019.
- 351 John Hewitt, Kawin Ethayarajh, Percy Liang, and Christopher D. Manning. Conditional probing: Measuring  
352 usable information beyond a single probe. In *Proceedings of NAACL*, 2021.
- 353 John Hopfield. Neural networks and physical systems with emergent collective computational abilities.  
354 *Proceedings of the National Academy of Sciences*, 1982.
- 355 Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. What does bert learn about the structure of language?  
356 In *Proceedings of ACL*, 2019.
- 357 Thomas Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *Inter-*  
358 *national Conference on Learning Representations*, 2017.
- 359 Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection  
360 for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 361 Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of  
362 words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*,  
363 2013a.
- 364 Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word repre-  
365 sentations. In *Proceedings of NAACL*, 2013b.
- 366 Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representa-  
367 tion. In *Empirical Methods in Natural Language Processing*, 2014.

- 368 Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexan-  
 369 der Miller. Language models as knowledge bases? In *Empirical Methods in Natural Language Processing*,  
 370 2019.
- 371 Tiago Pimentel, Josef Valvoda, Rowan Hall Maudslay, Ran Zmigrod, Adina Williams, and Ryan Cotterell.  
 372 Information-theoretic probing for linguistic structure. In *Proceedings of ACL*, 2020.
- 373 Abhilasha Ravichander, Yonatan Belinkov, and Eduard Hovy. Probing the probing paradigm: Does probing  
 374 accuracy entail task relevance? In *Proceedings of EACL*, 2021.
- 375 Lucrezia Valeriani, Diego Doimo, Francesca Cuturello, Alessandro Laio, Alessio Ansuini, and Alberto  
 376 Cazzaniga. The geometry of hidden representations of large transformer models. *arXiv preprint*  
 377 *arXiv:2302.00294*, 2023. Presented as a poster at NeurIPS 2023.
- 378 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep  
 379 learning requires rethinking generalization. *International Conference on Learning Representations*, 2017.

## 380 8 Appendix

### 381 A Experimental Setup, Resources, and Reproducibility

382 This section documents the experimental setup, hyperparameters, computational resources, and implemen-  
 383 tation details required to reproduce the results in the paper. The goal is to ensure transparency and  
 384 reproducibility rather than to optimize for computational efficiency.

385 **Model architecture and size.** All experiments use a decoder-only Transformer trained from scratch.  
 386 The architecture follows a GPT-2-style decoder with 8 Transformer layers, each employing multi-head self-  
 387 attention with 4 attention heads over a  $d$ -dimensional embedding space. Embedding dimensionality is varied  
 388 across experiments ( $d \in \{32, 64, 128, 256\}$ ), and graph size is varied over  $|V| \in \{15, 60, 150\}$ . Each node in  
 389 the graph is mapped to a unique discrete token, and the model predicts successor tokens given a single-  
 390 node input. As a reference configuration, for  $d = 128$  and  $|V| = 60$  (the primary setting reported in the  
 391 main paper), the model contains approximately 1.5–1.7 million parameters, depending on bias terms and  
 392 embedding size.

393 **Training procedure and hyperparameters.** Models are trained using the AdamW optimizer with  
 394 learning rate  $1 \times 10^{-3}$ , weight decay 0.1, and batch size 32. The training objective is cross-entropy loss  
 395 over successor prediction. All runs are trained for a fixed 10,000 epochs, which is sufficient to reach stable  
 396 convergence and near-perfect set accuracy ( $\geq 99.9\%$ ) across all graph families. We do not employ early  
 397 stopping, validation-based model selection, or automated hyperparameter search. Hyperparameter values  
 398 were chosen based on standard practice and empirically verified to yield stable convergence across graph  
 399 families and model sizes. Hyperparameters were not tuned on held-out evaluation data.

400 **Graph construction and data generation.** Experiments are conducted on synthetic relational graph  
 401 worlds constructed using standard graph generators. For chain, cycle, star, and 3-regular graphs, successor  
 402 relations are directed. For Watts–Strogatz graphs, edges are undirected and represented using bidirectional  
 403 successor relations. Watts–Strogatz graphs are generated with nearest neighbors  $k = 4$  and rewiring prob-  
 404 ability  $p = 0.2$ . Unless otherwise stated, graphs contain  $|V| = 60$  nodes. Shortest-path distances used in  
 405 geometric probes are computed with respect to the graph definition for each topology, and node pairs that  
 406 are unreachable under the corresponding graph are excluded from distance-based analyses.

407 **Random seeds and reporting.** For the primary configuration reported in the main paper ( $d = 128$ ,  
 408  $|V| = 60$ ), all quantitative results are aggregated over 15 independent random seeds and reported as mean  
 409  $\pm$  standard deviation. For all other configurations (embedding-dimension sweeps, graph-size sweeps, and  
 410 ablations reported in the Appendix), results are reported from a single representative run. Unless otherwise  
 411 stated, all reported metrics correspond to the final Transformer layer after convergence.

412 **Computational resources.** All experiments were run on Google Colab CPU instances. Training time var-  
 413 ied substantially depending on graph topology, graph size, embedding dimension, and random seed. Even the  
 414 smallest configurations required on the order of an hour per training run, while larger or higher-dimensional  
 415 settings required multiple hours. The full 15-seed sweep for the  $d = 128$ ,  $|V| = 60$  configuration required  
 416 approximately two days of computation. Across all experiments, including embedding-dimension sweeps,  
 417 graph-size sweeps, loss-function ablations, and multiple random seeds, the total computational budget is  
 418 estimated at several days of CPU computation.

419 **Software and implementation.** Models are implemented in PyTorch using standard Transformer com-  
 420 ponents. Graph construction and shortest-path computations use NetworkX. UMAP is used exclusively for  
 421 qualitative visualization and is not used for quantitative evaluation or for any claims in the paper. All  
 422 experiments rely on publicly available libraries, and no external codebases were modified. Code to reproduce  
 423 all experiments will be made available upon acceptance.

## 424 B Additional Experimental Results

425 This appendix provides supplementary experimental results supporting the claims in the main paper. Unless  
 426 otherwise stated, all values are reported from the *final layer* of the model (Layer 8), reflecting representations  
 427 at convergence.

428 All tables use the same geometric probes as in the main paper. Higher Spearman correlation and  $R^2$   
 429 indicate stronger geometric alignment with graph shortest-path distance, while lower distortion indicates  
 430 greater metric uniformity. Expansion variance and local isometry provide complementary diagnostics of  
 431 low-dimensional structure and neighborhood preservation.

### 432 B.1 Embedding Dimension Sweeps (Cross-Entropy Loss)

433 We examine how embedding dimensionality affects geometric organization when memorization accuracy is  
 434 held constant. Graph size is fixed at  $n = 60$ , and all models are trained using cross-entropy loss. Across all  
 435 settings, models achieve perfect memorization, allowing us to isolate representational geometry from task  
 436 performance.

437 **Low-dimensional regime ( $d = 32$ ).** Results are reported in Table 2. At low dimensionality, represen-  
 438 tational capacity is limited, restricting the set of embedding configurations that can achieve perfect memo-  
 439 rization. In Chain and Cycle graphs, this constraint does not induce geometry the successor-prediction task  
 440 exposes no global ordering, and nodes remain largely interchangeable. As a result, embeddings fail to reflect  
 441 graph distances despite the apparent simplicity of these topologies.

442 In contrast, the 3-regular graph already exhibits measurable geometric alignment at  $d = 32$ . Although  
 443 all nodes have identical degree, higher-order connectivity patterns differ, introducing weak but genuine  
 444 asymmetry. Under limited capacity, these differences partially constrain representations, yielding non-zero  
 445 rank correlation. However, scaling remains inconsistent, indicating that this alignment is local rather than  
 446 globally metric.

447 Watts–Strogatz graphs show clearer geometric structure even at low dimensionality. Long-range shortcuts  
 448 introduce global asymmetry by differentiating nodes according to path centrality and multi-hop reachability,  
 449 allowing embeddings to reflect graph distances even under tight capacity constraints.

450 **Intermediate dimensionality ( $d = 64$ ).** Results are reported in Table 3. As dimensionality increases,  
 451 representational capacity expands. For Chain and Cycle graphs, this does not improve geometry: the  
 452 task still permits index-based solutions, and additional dimensions only increase the space of non-metric  
 453 representations.

454 For the 3-regular graph, increased capacity weakens geometric alignment relative to  $d = 32$ . Though local  
 455 neighborhoods remain structured, the graph’s near-regularity permits many nodes to be represented by

456 shared functional roles rather than arranged according to global graph distance. The model can satisfy the  
 457 task using these alternative encodings, leading to reduced rank consistency.

458 Watts–Strogatz graphs, by contrast, benefit from increased dimensionality at this stage. The presence of  
 459 shortcuts continues to break symmetry at multiple scales, and additional capacity supports a more stable  
 460 embedding that better reflects global graph distances.

461 **High-dimensional regime ( $d = 256$ ).** Results are reported in Table 4. At high dimensionality, these  
 462 trends become more pronounced. Chain and Cycle graphs remain non-geometric, as additional capacity  
 463 further decouples representations from any implicit ordering.

464 In the 3-regular graph, geometric alignment deteriorates further. Without strong global anchors, the model  
 465 increasingly exploits high-dimensional freedom to realize symmetry-equivalent memorization strategies. This  
 466 results in unstable geometry characterized by reduced rank alignment, increased distortion, and noticeable  
 467 variability across random seeds.

468 Watts–Strogatz graphs, however, maintain stable geometric organization at higher dimensions. Once  
 469 shortcut-induced asymmetry establishes a meaningful global structure, additional dimensions neither dis-  
 470 rupt nor substantially enhance geometry. Instead, embeddings saturate into a consistent configuration that  
 471 preserves graph distances under our probes.

472 **Summary.** Embedding dimensionality interacts with graph topology in a graph-specific manner. Highly  
 473 interchangeable structures such as Chains and Cycles remain non-geometric across all dimensions. Graphs  
 474 with weak but higher-order structural asymmetry, such as 3-regular graphs, exhibit dimension-sensitive and  
 475 seed-unstable geometry. In contrast, graphs with persistent multi-scale asymmetry, such as Watts–Strogatz  
 476 graphs, support stable geometric organization across a wide range of embedding dimensions. The star graph  
 477 forms a distinct boundary case its degenerate hub–spoke distance structure imposes a trivial geometric  
 478 separation that remains unchanged across embedding dimensions. These results highlight that changes in  
 479 embedding dimension affect geometry through their interaction with specific structural features, rather than  
 480 through capacity alone.

## 481 B.2 Loss Function Ablation

482 We next test whether representational collapse is specific to the cross-entropy objective. We compare against  
 483 a uniform  $1/K$  loss that assigns equal probability mass to all valid successors, holding graph size ( $n = 60$ )  
 484 and embedding dimension ( $d = 128$ ) fixed. All comparisons are reported for the *final layer* only, after models  
 485 achieve stable convergence and near-perfect set accuracy. Aggregate geometric probe results are reported in  
 486 Table 5.

487 Across all graph families, the choice of loss function does not qualitatively alter whether geometric structure  
 488 emerges at convergence. Graph families with high node interchangeability (Chain and Cycle) remain non-  
 489 geometric under both objectives, exhibiting near-zero rank correlation and extremely high distortion despite  
 490 perfect memorization. This indicates that representational collapse in these settings is not an artifact of  
 491 the cross-entropy loss, but rather a consequence of the task admitting multiple functionally equivalent,  
 492 non-metric solutions.

493 For the 3-regular graph, which introduces partial asymmetry while maintaining uniform degree, the uniform  
 494  $1/K$  objective does not reliably induce recoverable geometry. In the reported run, rank correlation remains  
 495 modest and global scaling inconsistency persists, with moderate distortion and high expansion variance.  
 496 Given the weak and unstable geometry already observed under cross-entropy, these results indicate that  
 497 partial symmetry alone is insufficient to support a coherent global metric embedding across objectives, and  
 498 that apparent differences may reflect seed-level variability rather than a systematic effect of the loss function.

499 In the Watts–Strogatz graph, which strongly breaks global symmetry while preserving structured local  
 500 connectivity, robust geometric organization emerges under both loss functions. Final-layer rank correlation  
 501 and  $R^2$  remain high, and distortion remains controlled, demonstrating that sufficiently strong relational

constraints induce geometric structure regardless of the training objective. The  $1/K$  loss primarily affects the smoothness of output distributions rather than the presence of geometry itself.

Finally, the star graph exhibits degenerate geometric behavior under both objectives. Rank-based metrics are undefined due to the lack of variance in graph distances, while distortion and expansion variance reflect topology-imposed constraints rather than learned metric organization. Under the uniform  $1/K$  objective, local isometry is higher, reflecting the trivial and identical 1-hop neighborhoods imposed by the hub–spoke topology. This does not indicate meaningful geometric structure, but rather highlights that purely local probes become ill-conditioned in settings with extreme distance degeneracy.

**Summary.** Altering the training objective may modulate the stability of weak geometric effects, but does not determine whether geometric structure emerges. Recoverable geometry at convergence is governed primarily by the topological and symmetry constraints of the underlying graph, rather than by the choice of loss function.

### B.3 Effect of Graph Size

Finally, we evaluate whether representational collapse is a finite-size artifact by varying graph size under cross-entropy training.

**Small graphs.** Results for smaller graphs are reported in Table 6. At small graph sizes ( $n = 15$ ), geometric alignment emerges selectively depending on the degree of relational identifiability imposed by graph topology. Graph families with limited node interchangeability exhibit measurable geometric alignment, while graph families with high interchangeability continue to admit non-geometric memorization solutions (Table 6). In this small- $n$  regime, each relational constraint influences a large fraction of node pairs, making global geometric probes particularly sensitive to weak symmetry-breaking effects.

At small  $n$ , partial asymmetries introduced by graph topology constrain the space of admissible representations more strongly, as fewer functionally equivalent embeddings satisfy the training objective. This finite identifiability regime amplifies weak topological constraints that would otherwise be diluted at larger scales. As a result, even incomplete relational structure can induce measurable rank alignment between embedding distances and graph distances.

This effect is most evident in the 3-regular graph. Although all nodes share identical degree, higher-order connectivity patterns differ substantially at  $n = 15$ , including variations in cycle participation, neighborhood overlap, and multi-hop reachability. These differences reduce node interchangeability and limit the viability of purely index-based memorization strategies. Consequently, the model exhibits relatively high Spearman rank correlation, indicating partial geometric alignment. However this geometry remains fragile, distortion and expansion variance remain elevated and not topology-consistent, reflecting inconsistent global scaling and the absence of a coherent metric embedding despite non-zero rank alignment.

In contrast, the Watts–Strogatz graph does not exhibit strong geometric alignment at  $n = 15$ . While the topology introduces long-range shortcuts that break global symmetry, at small  $n$  these shortcuts are too sparse to establish a stable global ordering of nodes. Local ring structure dominates, and many nodes remain functionally similar under the successor prediction objective. As a result, the model can satisfy the task using localized memorization strategies that do not require embedding distances to reflect global graph distances. This explains the low rank correlation observed despite relatively bounded distortion, indicating that low distortion alone is insufficient to guarantee meaningful geometric organization.

Highly symmetric graph families such as chains and cycles continue to admit non-geometric memorization strategies even at small scale. In these settings, nodes remain functionally interchangeable under the training objective, allowing the model to achieve perfect memorization without encoding relational distances. Accordingly, rank correlation remains near zero and distortion remains high, indicating representational collapse rather than metric structure.

The star graph represents a distinct boundary case. At  $n = 15$ , graph-theoretic distances are degenerate, consisting primarily of hub–spoke relations with no meaningful ordering among spoke nodes. Ratio-based

549 geometric probes such as distortion are therefore ill-conditioned in this regime individual spoke placements  
 550 exert disproportionate influence on extreme distance ratios, leading to elevated distortion values. Impor-  
 551 tantly, this does not reflect a failure to learn task-relevant structure. The model learns a task-aligned  
 552 geometric separation between the hub and the set of spokes that is sufficient for perfect successor prediction,  
 553 while distances among spokes remain unconstrained and unidentifiable.

554 Together, these observations demonstrate that geometric structure does not arise from small scale alone.  
 555 Rather, small- $n$  regimes amplify the effect of topological constraints, revealing partial geometry when rela-  
 556 tional asymmetries restrict the solution space, and exposing metric degeneracies when relational structure  
 557 remains underdetermined.

558 **Large graphs.** Results for larger graphs are reported in Table 7. At larger graph sizes ( $n = 150$ ), the  
 559 behavior of geometric probes diverges sharply across graph families, reflecting how each topology interacts  
 560 with scale under the memorization objective.

561 For the 3-regular graph, the partial geometric alignment observed at smaller  $n$  disappears entirely. Although  
 562 node degree remains fixed, increasing graph size renders the topology progressively more homogeneous most  
 563 nodes occupy statistically similar positions with respect to local neighborhoods and multi-hop connectivity.  
 564 As  $n$  grows, these similarities dominate the task signal, making many nodes functionally indistinguishable  
 565 from the model’s perspective. The successor prediction objective can therefore be satisfied by grouping nodes  
 566 into a small number of role-based equivalence classes rather than organizing them in a globally consistent  
 567 metric space. This leads to representational collapse, characterized by near-zero rank correlation and elevated  
 568 distortion despite perfect memorization.

569 Chain and cycle graphs exhibit a different failure mode. While these topologies possess a well-defined global  
 570 metric structure, the successor prediction task does not require the model to preserve long-range distances  
 571 as graph size increases. As  $n$  grows, the majority of training signal is dominated by local transitions, and  
 572 errors in global ordering become increasingly irrelevant to the loss. Consequently, embeddings need only  
 573 preserve local adjacency rather than coherent global geometry, resulting in vanishing rank correlation and  
 574 severe distortion at scale.

575 In contrast, Watts–Strogatz graphs maintain stable geometric alignment even at large  $n$ . Although lo-  
 576 cally near-regular, this topology introduces persistent structural heterogeneity through long-range shortcuts.  
 577 As the graph grows, differences in shortcut participation, features prevent large equivalence classes of inter-  
 578 changeable nodes from forming, forcing the model to preserve global distance relationships in order to satisfy  
 579 the training objective. As a result, geometric alignment remains stable, with moderate rank correlation and  
 580 controlled distortion.

581 Finally, star graphs represent a degenerate boundary case. The presence of a single hub imposes an explicit  
 582 and unavoidable geometric separation between the center and the leaves. This structure is trivially imposed  
 583 by the topology and does not depend on learning or scale. After excluding unreachable node pairs, star  
 584 graphs admit only a single non-trivial distance class, rendering rank-based geometric probes degenerate by  
 585 construction. Accordingly, star graphs are excluded from correlation-based analysis.

586 Together, these results show that increasing graph size does not inherently promote geometric organization.  
 587 Instead, scale interacts with topology by either diluting or preserving the structural signals required for  
 588 geometric identifiability. Topologies whose relational constraints remain informative as  $n$  grows continue to  
 589 support recoverable geometry, while those that become functionally homogeneous or locally sufficient under  
 590 the task objective admit non-geometric memorization solutions.

591 **Summary.** Across graph families and scales, increasing graph size alone does not induce geometric orga-  
 592 nization. When topological constraints are weak or homogeneous, perfect memorization is possible without  
 593 distance-aware representations, leading to probe collapse. In contrast, topologies with persistent structural  
 594 heterogeneity support recoverable geometry. These results show that geometric structure is governed by  
 595 topological identifiability rather than graph size or model capacity.

Graph	Spearman	$R^2$	Distortion	Exp. Var.	Local Iso.
CHAIN	0.0657	0.0017	58.67	8.2647	14.3473
CYCLE	0.1342	0.0164	65.02	7.0019	13.2222
3-REG	0.5257	0.1569	7.20	16.1035	12.8129
WS	0.5685	0.3748	7.13	4.5539	11.2405
STAR	–	–	1.13	0.0000	0.0985

Table 2: Geometry diagnostics for  $n = 60$ ,  $d = 32$  using the final layer (Layer 8).

#### B.4 Qualitative Visualization of Learned Representations

To complement the quantitative geometric probes reported above, we provide qualitative visualizations of learned representations for each graph family. These visualizations are intended as illustrative diagnostics rather than primary evidence, as low-dimensional projections may distort global geometry. Accordingly, all claims in the paper are grounded in the distance-based probes reported in the main text and appendix tables.

Each dashboard visualizes (i) a 2D projection of node embeddings, (ii) a 3D embedding visualization, (iii) the underlying relational graph for a representative node, and (iv) the corresponding softmax successor probability distribution. All visualizations are generated from graphs with 60 nodes, using embeddings of dimension 128.

## C Architecture Independence: Encoder-Decoder Transformers

In the main text, our analyses focus on decoder-only Transformers to align with common architectures used in autoregressive language modeling. To verify that the observed relationship between memorization, topological symmetry, and geometric emergence is a fundamental property of the learning objective rather than an architectural artifact, we replicate our core experiments using an encoder-decoder framework. We adopt a standard encoder-decoder Transformer trained from scratch under the identical successor-prediction objective. To maintain capacity comparability with the decoder-only experiments—which utilized 8 layers—we configure the encoder-decoder model to have 4 encoder layers and 4 decoder layers, yielding an equivalent total depth. The embedding dimension is fixed to  $d = 128$  with 4 attention heads, and the graph size to  $n = 60$ . As with the primarily analyzed models, the encoder-decoder configurations achieve near-perfect memorization (set accuracy  $\geq 99.9\%$ ) across all tested graph families.

We probe the internal representations to determine whether introducing an encoder-decoder architecture alters geometric recoverability. Our findings demonstrate that the representations exhibit the precise pattern of geometry emergence observed in decoder-only models (Table 8). In highly symmetric topologies, such as chains and cycles, the model does not develop a coherent global metric. The learned embeddings exhibit near-zero rank correlation with graph shortest-path distances and high distortion, confirming that high node interchangeability permits non-geometric memorization regardless of the encoder-decoder architecture. Under partial symmetry-breaking, such as in 3-regular graphs, the representation space shows only weak and fragile geometric alignment, characterized by inconsistent global scaling.

Conversely, under the strong symmetry-breaking of Watts–Strogatz small-world graphs, coherent and stable geometric structure robustly emerges within the model’s internal representations. In the boundary case of star graphs, the geometry degenerates into the trivial, topology-dictated hub-and-spoke separation, confirming it as an imposed limit rather than a learned spatial organization. These results establish that the selective emergence of geometric structure is architecture-agnostic. The non-identifiability of metric structure in symmetric graphs is a direct mathematical consequence of how the task objective interacts with relational topology, rather than an artifact of causal masking or decoder-only processing. Consequently, the conclusion that geometric structure depends primarily on symmetry constraints applies broadly across Transformer architectures.

Graph	Spearman	$R^2$	Distortion	Exp. Var.	Local Iso.
CHAIN	0.0105	0.0000	88.39	10.4470	16.2047
CYCLE	0.0413	0.0001	60.95	11.2844	16.9093
3-REG	0.3585	0.0847	7.69	21.6430	14.1518
WS	0.6077	0.4043	8.51	5.9267	12.8957
STAR	–	–	1.14	0.0000	0.0197

Table 3: Geometry diagnostics for  $n = 60$ ,  $d = 64$  using the final layer (Layer 8).

Graph	Spearman	$R^2$	Distortion	Exp. Var.	Local Iso.
CHAIN	0.0324	0.0006	170.98	22.3411	23.7879
CYCLE	0.0339	0.0012	94.06	26.9706	26.3165
3-REG	0.2052	0.0099	16.93	22.2633	15.0140
WS	0.6682	0.4634	8.17	5.4428	12.3224
STAR	–	–	2.65	0.0000	0.0192

Table 4: Geometry diagnostics for  $n = 60$ ,  $d = 256$  using the final layer (Layer 8).

Graph	Spearman	$R^2$	Distortion	Exp. Var.	Local Iso.
CHAIN	-0.0361	0.0017	83.13	15.6146	20.0585
CYCLE	0.0596	0.0032	66.41	14.7997	19.6616
3-REG	0.2740	0.0378	15.05	24.2272	15.8122
WS	0.4994	0.3204	9.01	8.7537	15.4345
STAR	–	–	1.05	0.000	0.9402

Table 5: Geometry diagnostics using uniform  $1/K$  loss ( $n = 60$ ,  $d = 128$ ), final layer (Layer 8).

Graph	Spearman	$R^2$	Distortion	Exp. Var.	Local Iso.
CHAIN	-0.2126	0.0042	16.70	12.7608	12.0401
CYCLE	0.1841	0.0477	14.96	10.3763	11.0858
3-REG	0.5590	0.1436	9.42	10.4900	8.4305
WS	0.0563	0.0032	4.96	9.9701	10.6380
STAR	–	–	16.01	1.1222	0.6594

Table 6: Geometry diagnostics for small graphs ( $n = 15$ ,  $d = 32$ ), final layer (Layer 8).

Graph	Spearman	$R^2$	Distortion	Exp. Var.	Local Iso.
CHAIN	0.0053	0.0000	201.75	5.6996	17.2706
CYCLE	0.0384	0.0014	201.44	6.2304	18.1359
3-REG	-0.0265	0.0005	28.03	23.8432	15.1244
WS	0.5383	0.3879	10.51	5.3564	16.1860
STAR	–	–	1.25	0.0000	0.0046

Table 7: Geometry diagnostics for large graphs ( $n = 150$ ,  $d = 384$ ), final layer (Layer 8).

Graph	Spearman	$R^2$	Distortion	Exp. Var.	Local Iso.
CHAIN	0.0184	0.0001	80.68	74.6450	44.8407
CYCLE	0.0066	0.0000	81.20	75.3669	45.2177
3-REG	0.2006	0.0075	6.67	142.4757	40.6776
WS	0.5950	0.4106	8.97	64.9011	41.4645
STAR	–	–	320.75	8.2636	6.0476

Table 8: Geometry diagnostics for encoder-decoder models ( $n = 60$ ,  $d = 128$ ), final layer (Decoder Layer 4).

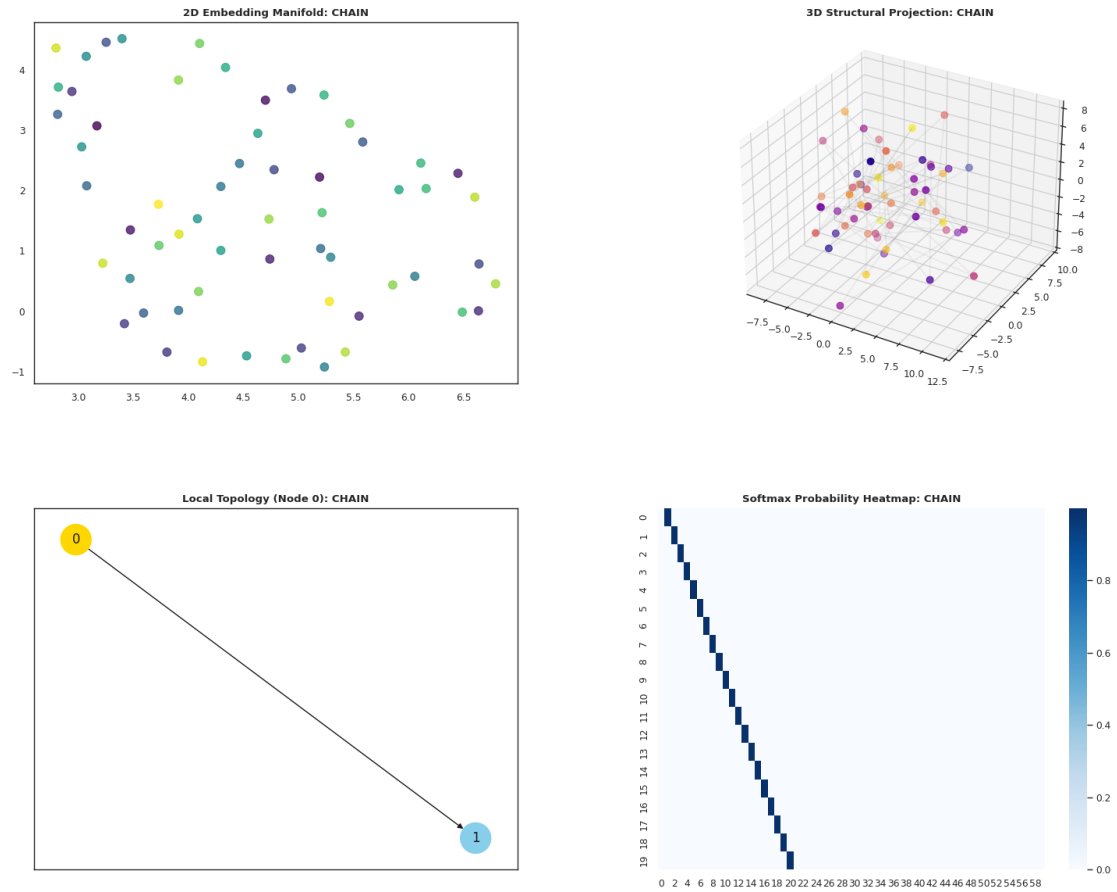


Figure 2: **Chain graph: memorization without geometric structure.** Embeddings do not exhibit a coherent global metric despite perfect task performance.

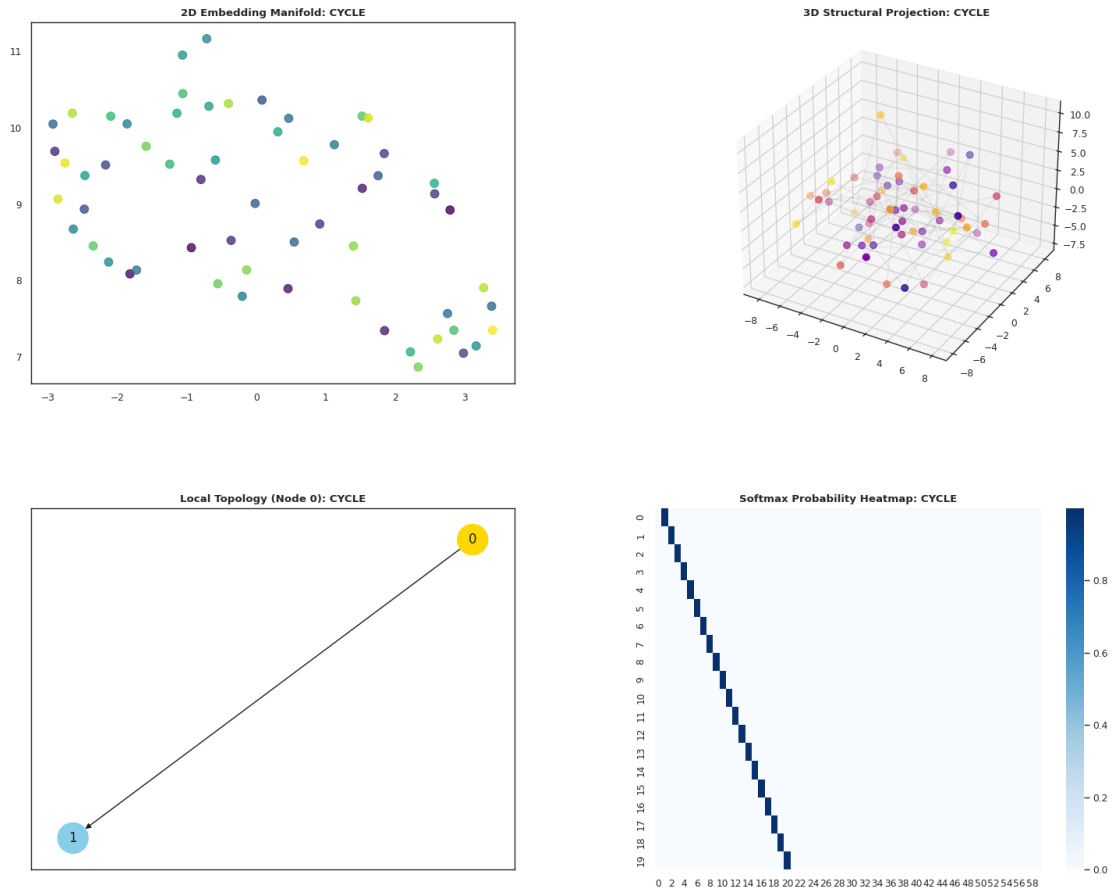


Figure 3: **Cycle graph: symmetry-induced non-geometric memorization.** Despite perfect memorization, embeddings do not exhibit a coherent global metric structure.

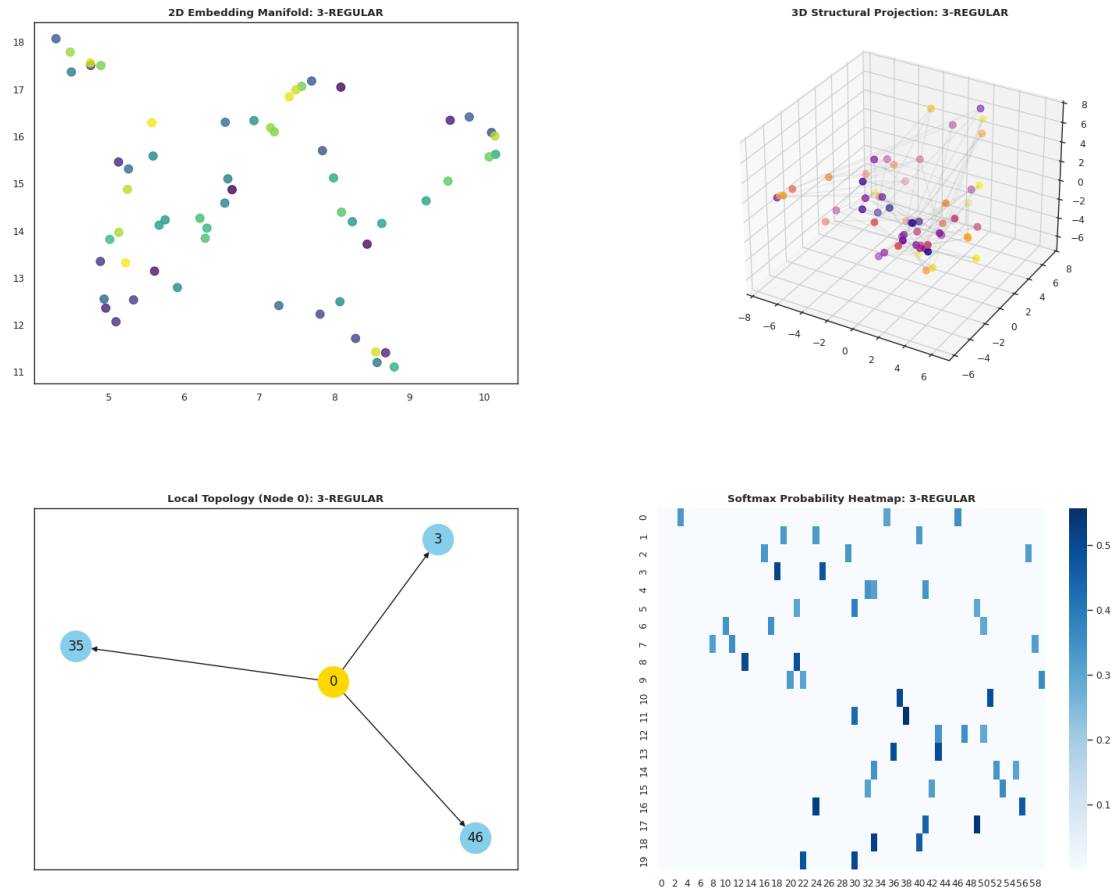


Figure 4: **3-Regular graph: partial asymmetry and fragile geometry.** Local neighborhood coherence emerges without a stable global metric, consistent with low-moderate rank correlation and elevated distortion.

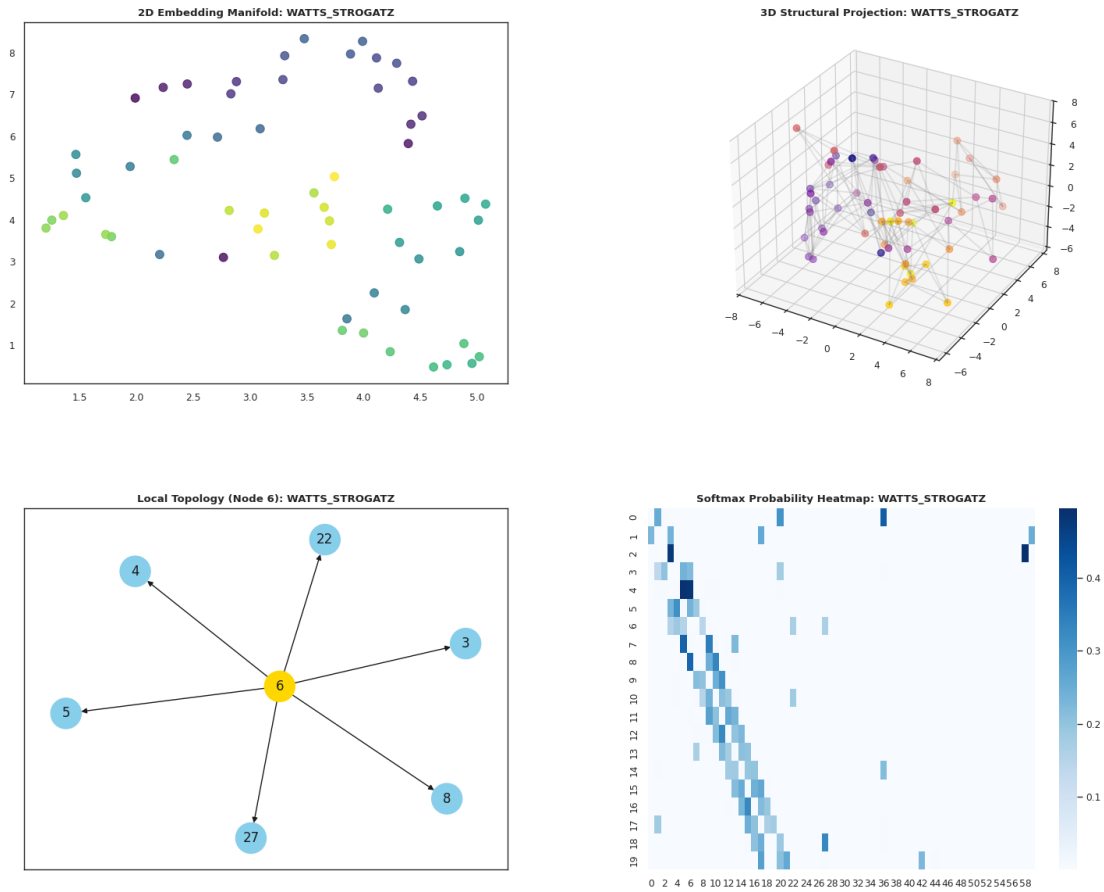


Figure 5: **Watts–Strogatz graph: emergent geometric organization.** Embeddings form a smooth low-dimensional structure that preserves graph distances, consistent with high rank correlation and controlled distortion.

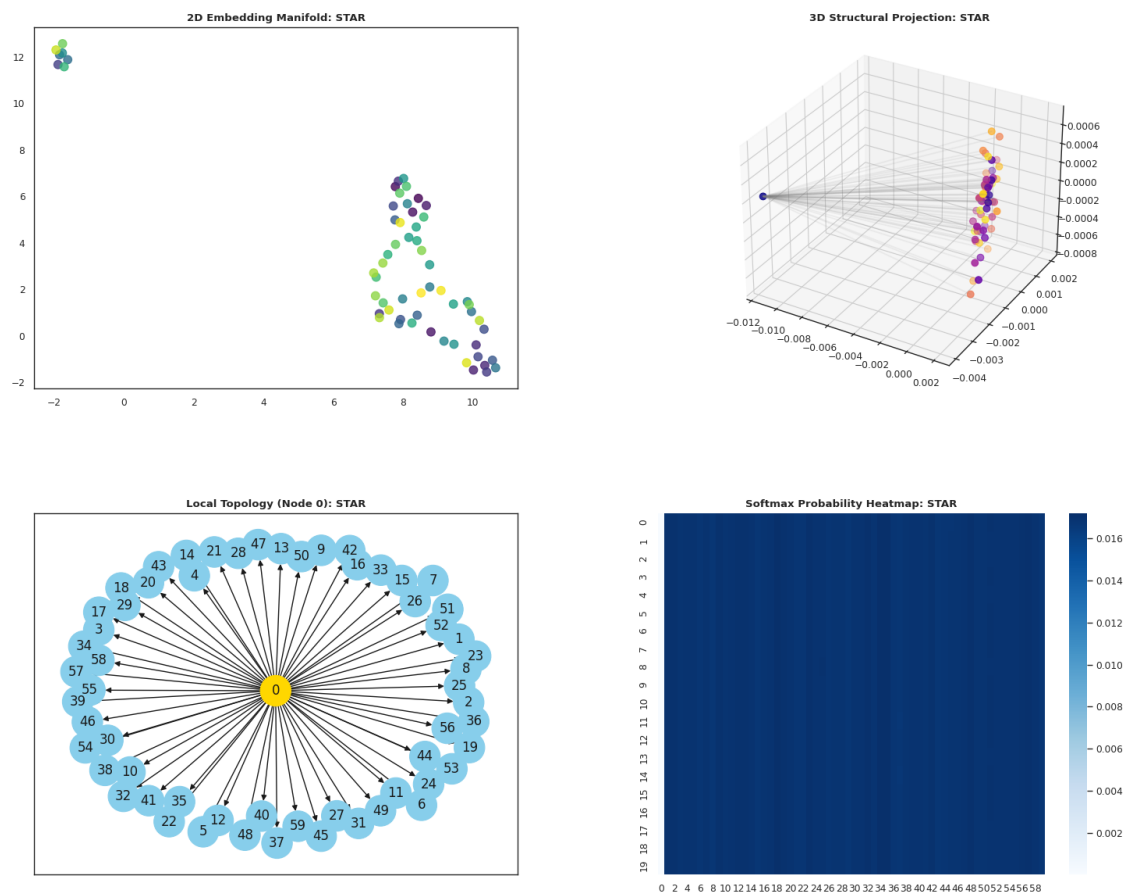


Figure 6: **Star graph: topology-determined geometric structure.** The hub–spoke topology yields a degenerate but task-sufficient geometric separation. Rank-based probes are undefined due to distance collapse, and the observed structure reflects topological degeneracy rather than recoverable metric organization.