

---

# The Connectome of a Large Language Model

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 A central question about information processing systems is whether fundamentally  
2 different optimization processes lead to similar internal organization. Cortical  
3 networks exhibit small-world topology, modular communities, scale-free hubs,  
4 and structural-functional dissociation, and sparse transcoders now provide the  
5 interpretable units needed to test these same properties inside a large language  
6 model. We use skip transcoders from Gemma Scope 2 to extract 16,384 sparse  
7 interpretable features per MLP layer of Gemma 3 (270M and 1B), construct per-  
8 layer co-activation networks from 1M tokens of FineWeb-Edu, and analyze them  
9 with a standard network-neuroscience toolbox augmented by the Laplacian renor-  
10 malization group and the spectral heat capacity. Both models exhibit analogues of  
11 all four classical cortical signatures, though with threshold sensitivity in the small-  
12 world case and quantitative differences in fractal dimension. Multi-scale analysis  
13 further reveals a dual-regime fractal scaling with a local exponent  $d_B \approx 1.35$   
14 below the Leiden community diameter, a single-peaked heat capacity that sharpens  
15 with depth, and a non-monotonic cross-layer coupling trajectory. Hub features  
16 and coarse-grained communities carry interpretable semantics organized along a  
17 tokenization-syntax-semantics-morphology hierarchy. These findings indicate that  
18 transformer representations trained on next-token prediction exhibit several network  
19 properties that are statistically similar to those observed in cortical networks.

## 20 1 Introduction

21 A central question for any field that studies complex information processing systems is whether  
22 systems shaped by very different processes converge on a common internal architecture. Biological  
23 brains have been shaped by hundreds of millions of years of evolution under tight energy and  
24 wiring constraints, and a long tradition in network neuroscience has shown that the resulting cortical  
25 networks exhibit a small set of robust organisational principles. These include small-world topology  
26 with high clustering and short path length [41, 4], modular community structure that segregates  
27 and integrates information [35, 27], scale-free hubs that dominate global communication [1], and  
28 a partial dissociation between structural and functional connectivity in which weight-based and  
29 correlation-based connectomes are only loosely aligned [19]. Large language models, in contrast, are  
30 shaped by gradient descent on text corpora over a few months of compute and have very different  
31 inductive biases. Whether the trained models nevertheless converge on the same network signatures  
32 is a question with theoretical motivation — modularity can emerge spontaneously under modularly  
33 varying objectives [14], and sparse network topologies inspired by network science have already  
34 improved ANN training [25] — but one that, until recently, could not be posed because the units of  
35 analysis inside an LLM were not well defined.

36 Existing approaches to LLM internals have focused on tracing causal circuits behind specific com-  
37 putations or comparing representational geometry with biological systems, but these approaches  
38 do not characterize the network organization of the model as a whole. Recent progress on sparse

39 decompositions of internal activations changes this picture. Sparse autoencoders and transcoders  
40 learn to reconstruct each MLP block from a much wider dictionary of features that activate sparsely  
41 on natural text [8, 38, 10, 29]. With the public release of Gemma Scope and Gemma Scope 2 [18, 24],  
42 every layer of the Gemma 3 family [36] now comes equipped with a per-layer skip transcoder whose  
43 features are largely interpretable in isolation. These features supply a candidate analogue of the  
44 neuron for network analysis. They are sparse, they fire on identifiable inputs, and they admit pairwise  
45 correlation analysis in the same way that neural spike trains do. We can therefore treat an LLM as an  
46 empirical network and apply the standard tools of network science to it, though the resulting networks  
47 depend on the choice of decomposition.

48 This paper takes that step. For two scales of the Gemma 3 family (270M and 1B), we build per-  
49 layer co-activation networks from transcoder features on a 1M-token corpus alongside weight-based  
50 structural networks from the transcoder parameters, and analyze both with the standard topology  
51 toolbox of network neuroscience augmented by multi-scale tools from statistical physics. We ask  
52 which of the classical cortical signatures hold, which fail, and where new structure appears. We then  
53 attach semantics to the topology by labelling hub features and coarse-grained communities through  
54 their maximally activating contexts.

55 The networks exhibit analogues of all four classical cortical signatures, and multi-scale analysis  
56 reveals three further properties, including a dual-regime fractal scaling with a local exponent of about  
57 1.35 below the Leiden community diameter, a single-peaked spectral heat capacity that sharpens  
58 with depth, and a non-monotonic cross-layer coupling trajectory. The quantitative differences from  
59 cortex are themselves informative. The LLM fractal dimension falls below the cortical value of  
60 approximately 2.0 reported by Gallos et al. [12], and the heat capacity shows a single sharp peak  
61 where cortical connectome data shows a broad plateau [31], both consistent with sparser and more  
62 specialized internal connectivity. Hub features and coarse-grained communities carry interpretable  
63 semantics organized along a depth-dependent tokenization-syntax-semantics-morphology hierarchy.  
64 To our knowledge, this is the first systematic network-science analysis of an interpretable LLM  
65 feature space at corpus scale.

## 66 **2 Related Work**

### 67 **2.1 Mechanistic interpretability and circuit analysis**

68 Mechanistic interpretability has made rapid progress in identifying the internal computations of  
69 language models through causal analysis. Ameisen et al. [2] introduced circuit tracing, which  
70 constructs attribution graphs that track how individual features and attention heads contribute to a  
71 model’s output on a given prompt, and Lindsey et al. [21] used thousands of such graphs to catalog  
72 recurring computational motifs. Marks et al. [23] discovered task-specific sparse feature circuits  
73 through causal mediation, and Dunefsky et al. [10] showed that transcoders facilitate the extraction  
74 of interpretable causal graphs. Because each attribution graph is tied to a single prompt and a single  
75 task, the cost of interpretation scales with the number of behaviors under study. These methods are  
76 powerful for explaining individual computations, and the resulting circuits have revealed rich internal  
77 structure, but they do not aggregate into a characterization of how the model as a whole is organized.

### 78 **2.2 Brain-AI representational comparison**

79 A separate line of work compares the internal representations of biological and artificial neural  
80 networks. Representational Similarity Analysis [16] and related methods such as Centered Kernel  
81 Alignment [15] quantify the similarity structure of population responses, providing a common  
82 framework for comparing brains and models across architectures. This approach has shown that task-  
83 optimized deep networks predict representational geometry in visual cortex [42] and that transformer  
84 language models trained on next-word prediction explain a large fraction of neural variance in the  
85 human language network [32]. Recent extensions incorporate topological features of representation  
86 spaces beyond pairwise distances [20]. RSA and related methods measure whether two systems  
87 encode similar information, and have established detailed correspondences between model layers and  
88 brain regions. Representational geometry is a static description of encoding, and does not capture the  
89 interactions among units that emerge during processing, as the well-documented structure-function  
90 dissociation in cortical networks illustrates [19].

## 91 2.3 Graph construction from LLM features

92 Several recent works have constructed graphs from LLM features without applying network-science  
93 analysis to them. Deng et al. [9] built co-activation graphs from sparse feature activations and  
94 showed that connected components correspond to composable semantic modules. Balcells et al. [3]  
95 tracked features across adjacent layers using Pearson and Jaccard scores on a 10M-token Pile sample,  
96 identifying pass-through, gating, and disappearing features. Laptev et al. [17] traced cross-layer  
97 feature flow through cosine similarity of decoder weights. The circuit tracing programme at Anthropic  
98 [2, 21] produced thousands of attribution graphs at the prompt level and used manual supernode  
99 grouping that effectively performs community detection by hand. None of these works compute  
100 corpus-wide network statistics on the resulting graphs, and none frames the analysis in terms of  
101 the cortical comparison. Two concurrent works move closer to network science. Zheng et al. [43]  
102 computes topological metrics on graphs built from raw neurons, and Liu et al. [22] applies ICA to  
103 neuron activations in the spirit of fMRI functional network analysis. Both operate on raw neurons  
104 without multi-scale topological analysis on interpretable features.

## 105 3 Data and Methods

106 We work with two scales from the Gemma 3 family, Gemma 3 270M (18 layers,  $d_{\text{model}} = 640$ ) and  
107 Gemma 3 1B (26 layers,  $d_{\text{model}} = 1152$ ), loading the corresponding Gemma Scope 2 skip transcoders  
108 with 16,384 features per layer [18, 24]. For each layer we build a functional network by streaming a  
109 1M-token sample of FineWeb-Edu [30] through the model, encoding MLP inputs with the transcoder,  
110 and accumulating a 16K by 16K Pearson co-activation matrix in a single pass. We complement  
111 these with structural networks built from the transcoder parameters, using cosine similarity between  
112 decoder weight vectors within a layer and the virtual weight  $M_{ij}^{(\ell \rightarrow \ell+1)} = |W_{\text{dec}}^{(\ell)}[i, :] \cdot W_{\text{enc}}^{(\ell+1)}[j, :]|$   
113 across layers [11, 10]. Both matrices are converted to sparse graphs by density-matched thresholding  
114 at a canonical density  $d = 0.001$ , with robustness checks at two additional densities. Details of the  
115 streaming computation, feature filtering, and threshold selection are in Appendix E.

116 The topology toolbox follows network neuroscience standards. We compute small-world indices,  
117 Leiden modularity at multiple resolutions [39], degree distribution statistics, assortativity, and the  
118 rich-club coefficient, testing every metric against both Erdős-Rényi and Configuration Model nulls  
119 [28]. Network robustness is assessed by percolation under random and targeted node removal [1],  
120 and structure-function coupling follows the methodology of Liégeois et al. [19], combining a Mantel  
121 correlation with a conditional analysis of functional connectivity at the strongest structural edges. For  
122 multi-scale structure we apply three renormalization group tools. Box covering [33, 34] yields the  
123 fractal dimension  $d_B$  and, through a piecewise fit, a local exponent below a crossover length. The  
124 Laplacian renormalization group of Villegas et al. [40] coarse-grains the network at a fixed diffusion  
125 time into supernodes, and the spectral heat capacity  $C(\tau) = \tau^2[\langle \lambda^2 \rangle_\tau - \langle \lambda \rangle_\tau^2]$  [31] diagnoses whether  
126 the network has a single characteristic scale, approximate scale invariance, or multiple well-separated  
127 scales. Full parameter settings and null model protocols are in Appendix E.

128 To attach semantics to the topology, we identify the top hub features by degree and the largest  
129 Laplacian renormalization group communities in each layer, and record their maximally activating  
130 contexts on a held-out corpus. Per-feature word clouds aggregate content words weighted by activation  
131 magnitude.

## 132 4 Results

133 We organize our results around five progressive questions about the internal organization of LLM  
134 feature networks.

### 135 4.1 Do feature co-activation networks have non-trivial structure?

136 At the canonical density  $d = 0.001$ , the Telesford small-world index  $\omega$  ranges over  $[-0.17, +0.57]$   
137 for Gemma 3 270M and  $[-0.43, +0.32]$  for Gemma 3 1B, with the layerwise mean close to zero [37].  
138 Early layers tend toward positive  $\omega$  (shorter paths relative to clustering) while deep layers shift toward  
139 negative  $\omega$  (higher clustering relative to path lengths), consistent with the increasing modularity  
140 reported in the next section. This property is threshold-sensitive and holds robustly only at sparse

141 densities ( $d \leq 0.002$ ). At  $d = 0.005$ ,  $\omega$  shifts substantially toward positive values in both models  
 142 (Appendix Figure 8a), indicating that the small-world balance is carried by the sparsest, strongest  
 143 connections.

144 Degree distributions are heavy-tailed across all layers, with skewness typically in the range 2 to 8 and  
 145 maximum degree reaching up to  $47\times$  the mean (270M) and  $53\times$  the mean (1B) in early layers. These  
 146 hub features have direct structural consequences for graph connectivity. Progressive node removal  
 147 reveals that networks remain largely intact under random failure (the largest connected component  
 148 retains approximately 40% of nodes after 50% removal) but collapse rapidly when high-degree or  
 149 high-betweenness nodes are targeted (Figure 2). This differential vulnerability to targeted versus  
 150 random attack is a defining signature of heterogeneous, hub-dependent architecture [1].

151 Null model comparisons confirm that this structure exceeds what degree heterogeneity alone can  
 152 explain. Configuration Model z-scores for modularity range from 198 to 635 in 270M and 68 to 808  
 153 in 1B, all far exceeding statistical significance. Feature co-activation networks are therefore highly  
 154 structured, with small-world topology, prominent hubs, and modular organization that cannot be  
 155 reduced to the degree sequence.

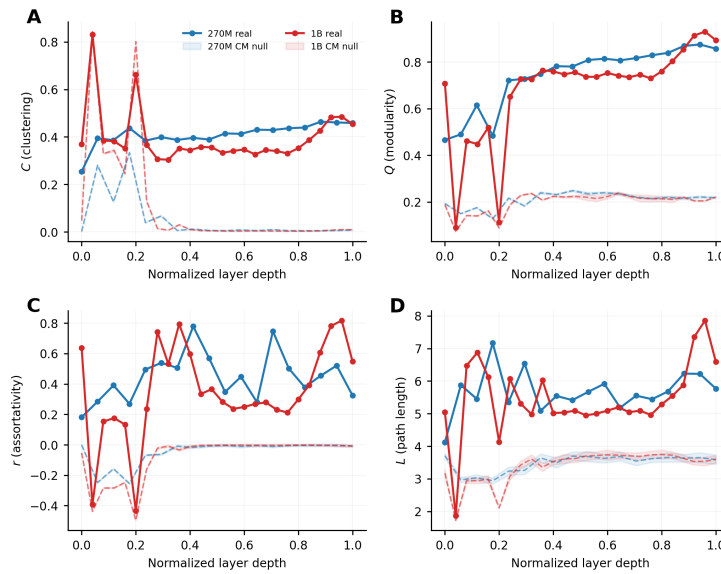


Figure 1: Layer evolution of four key topology metrics for Gemma 3 270M (blue) and 1B (red) at density-matched Pearson networks ( $d = 0.001$ ). Solid lines show observed values, dashed lines show the Configuration Model null mean with shaded  $\pm 2\sigma$  bands.

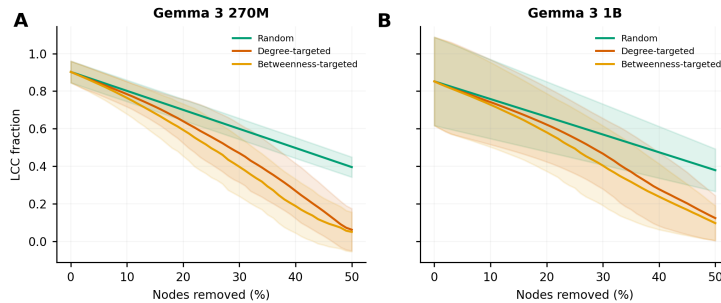


Figure 2: Network robustness under progressive node removal, averaged across all layers. **Left:** Gemma 3 270M. **Right:** Gemma 3 1B. Random removal (green) decays gradually, while degree-targeted (orange) and betweenness-targeted (red) removal collapse the largest connected component much faster. Shaded bands show  $\pm 1$  standard deviation across layers.

## 156 4.2 How does network organization evolve with layer depth?

157 Network organization could be uniform across layers or vary systematically with depth. Leiden  
158 modularity  $Q$  reveals a clear trend (Figure 1). In the 270M model,  $Q$  rises from 0.48 at layer 0 to  
159 0.88 at layer 16, a nearly twofold increase. In the 1B model,  $Q$  rises from 0.46 (layer 2) to 0.93  
160 (layer 24). Of all the properties we measure, this depth trend is the most consistent across density  
161 levels and both model scales. Configuration Model z-scores confirm that the trend reflects genuine  
162 organizational change, with modularity exceeding the degree-preserving null by hundreds of standard  
163 deviations at every layer (Appendix Table 2).

164 The monotonic increase in modularity implies that features in early layers form weakly differentiated  
165 groups with overlapping co-activation patterns, while features in deep layers form highly separated  
166 functional modules. This depth-dependent specialization is reminiscent of the cortical hierarchy  
167 in biological brains, where primary sensory areas show diffuse connectivity and association areas  
168 exhibit stronger modular organization [35].

169 Spectral analysis on the unthresholded Pearson matrices provides independent, threshold-free con-  
170 firmation. Fiedler values ( $\lambda_2$ ) increase from early to middle layers, reflecting growing algebraic  
171 connectivity (Appendix Figure 8b). These spectral signatures confirm that the structural evolution is  
172 a property of the correlation matrices themselves, independent of binarization.

173 Two layers in the 1B model (layers 1 and 5) are consistent outliers across multiple metrics, showing  
174 extremely low modularity ( $Q = 0.09$  and  $0.11$ ), high mean Pearson correlation, extreme hub degrees,  
175 and negative assortativity. These layers likely correspond to global attention layers in Gemma 3’s  
176 interleaved local/global attention architecture [36], where all features share the same context window  
177 and exhibit uniformly stronger co-activation. Their consistent anomalous behavior across every  
178 analysis serves as a useful internal control, suggesting that our metrics are sensitive to structural  
179 differences between layers.

## 180 4.3 What is the multi-scale geometry of these networks?

181 Biological neural networks exhibit multi-scale organization with fractal geometry within functional  
182 modules and small-world shortcuts between them [5, 12]. We ask whether LLM feature networks  
183 share this multi-scale organization.

184 Box-covering analysis of the thresholded graphs reveals two distinct scaling regimes in the relation  
185  $N_B(l_B) \sim l_B^{-d_B}$ . A single power-law fit yields  $d_B \approx 2.0$  with moderate  $R^2 \approx 0.93$  in both models,  
186 but allowing a piecewise fit with a breakpoint dramatically improves the fit and clears the  $\Delta\text{AIC} > 10$   
187 bar in most layers. Below the crossover length  $l_B^* \approx 6$ , the local fractal dimension averages  $1.37 \pm 0.14$   
188 (270M) and  $1.34 \pm 0.11$  (1B), consistent across models. Both models show a gradual increase in  $d_B$   
189 with depth, from approximately 1.0–1.25 in early layers to approximately 1.5 in late layers (Figure 3),  
190 indicating that deeper layers develop denser intra-module connectivity. Above the crossover, steeper  
191 scaling ( $d_B > 3$ ) reflects the small-world shortcuts connecting modules. The crossover length itself  
192 corresponds approximately to the diameter of Leiden communities, consistent with the fractal-module  
193 framework of Gallos et al. [12]. The local exponent of about 1.35 is lower than the cortical value  
194 of approximately 2.0, suggesting sparser internal connectivity inside LLM communities than inside  
195 cortical functional modules. With approximately ten data points per layer, the piecewise results are  
196 suggestive and the single-fit  $d_B$  with bootstrap CI serves as the primary conservative estimate.

197 Laplacian renormalization group coarse-graining probes compressibility at a fixed diffusion scale. A  
198 single round of LRG reduces the network to progressively fewer supernodes with increasing depth.  
199 Early layers yield 715–823 supernodes in 270M (9.7–11.2 $\times$  compression), while late layers yield as  
200 few as 191–209 supernodes (38–42 $\times$  compression). The 1B model shows even stronger late-layer  
201 compression, reaching up to 62.5 $\times$ . This trend, despite some layer-to-layer variability, is consistent  
202 with the increasing modularity reported above. Features within well-separated communities share  
203 diffusive modes and can be merged, so higher modularity produces greater compressibility. The  
204 supernode counts per depth are shown in Appendix Figure 10b.

205 The spectral heat capacity  $C(\tau)$ , computed from the full Laplacian spectrum, diagnoses whether  
206 the networks possess one characteristic scale, multiple well-separated scales, or approximate scale  
207 invariance [31]. All layers in both models show a single  $C(\tau)$  peak with no plateau, ruling out  
208 scale invariance in these networks. In the 270M model, the peak position  $\tau_{\text{peak}}$  decreases from 17.5

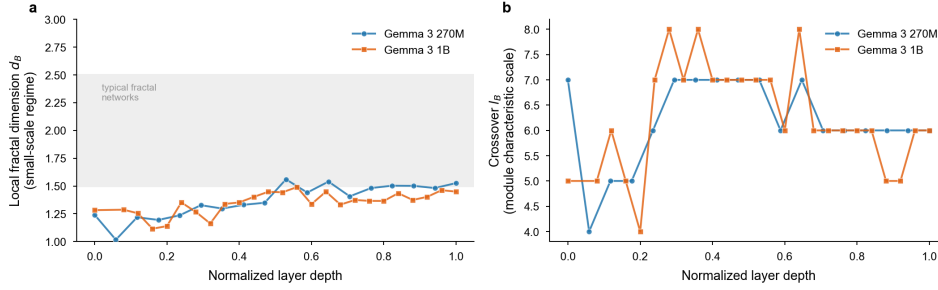


Figure 3: Box-covering fractal exponents across depth. **Left:** the local exponent  $d_B^{\text{local}}$  from the piecewise fit shows a gradual increase from approximately 1.0–1.25 in early layers to approximately 1.5 in late layers, with both models following the same trend. **Right:** the crossover length  $l_B^*$  has a median around 6, matching the empirical Leiden community diameter.

209 (layer 0) to 9.9 (layer 16), and the peak height  $C_{\text{peak}}$  increases from 4.8 to 8.9 across depth (Figure 4),  
 210 indicating that the dominant structural scale shifts to shorter diffusion times and grows in contrast. In  
 211 the 1B model,  $\tau_{\text{peak}}$  similarly decreases from 13.4 to 9.5. The  $C_{\text{peak}}$  values in the 1B model are more  
 212 variable, with the anomalous layers 0–1 and 4–5 showing unusually large values (12–16) and the  
 213 non-anomalous layers falling between 7.4 and 10.7 with a modest increase toward late layers. The  
 214 peak width (FWHM) narrows with depth in the 1B model, indicating convergence toward a single  
 215 well-defined scale, while the 270M model shows more variable FWHM across layers.

216 These three analyses converge on a coherent multi-scale picture. LLM feature networks are fractal  
 217 within modules ( $d_B \approx 1.35$ ) and connected by small-world shortcuts across modules. Each layer  
 218 possesses a single dominant structural scale, and this scale becomes sharper and more compact with  
 219 depth. The networks are modular and hierarchical, but they are not scale-invariant.

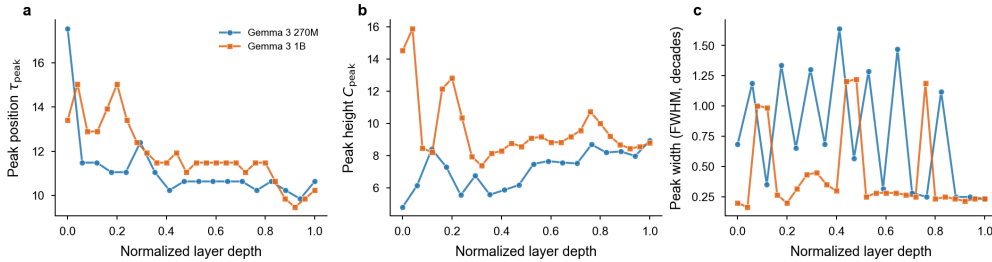


Figure 4: Evolution of spectral heat capacity peak parameters with depth. **Left:**  $\tau_{\text{peak}}$  decreases. **Center:**  $C_{\text{peak}}$  increases with depth in 270M, while 1B shows more variable values. **Right:** FWHM narrows in 1B but varies in 270M.

#### 220 4.4 How does information flow between layers?

221 Beyond within-layer organization, the residual stream mediates directed communication between  
 222 adjacent layers. Global correlation between structural (cosine) and functional (Pearson) connectivity  
 223 matrices is near zero across all layers (Figure 5), with Mantel  $r$  values ranging from  $-0.007$  to  
 224  $+0.013$  (270M) and  $-0.003$  to  $+0.034$  (1B). Despite this global dissociation, strong structural edges  
 225 carry a meaningful signal. The mean functional connectivity among the top 0.1% strongest structural  
 226 edges is approximately 20–25 times the global mean, and this enrichment decays monotonically  
 227 with relaxing threshold, reaching approximately  $2\times$  at the top 10%. The pattern parallels the well-  
 228 documented structure-function relationship in brain networks, where DTI structural connectivity and  
 229 fMRI functional connectivity show weak global correlation but strong local correspondence along  
 230 white-matter tracts [19]. A threshold-free Laplacian spectrum check with the same depth dependence  
 231 is reported in Appendix B.

232 Turning to directed coupling between adjacent layers, we quantify the cross-layer weight overlap  
 233 using `concentration_80`, the fraction of downstream LRG communities accounting for 80% of

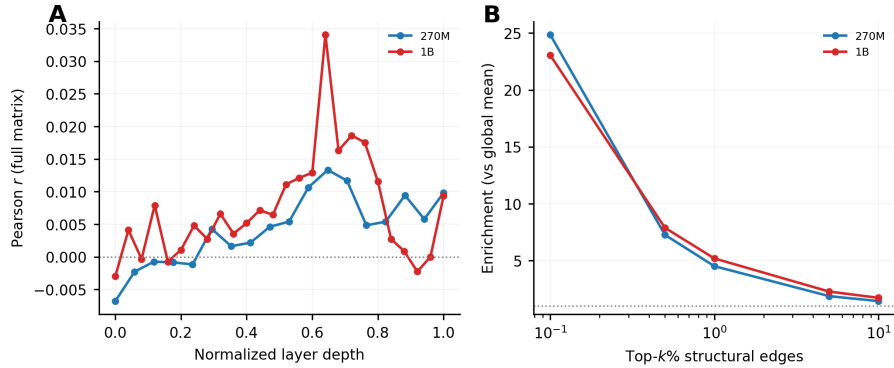


Figure 5: Structure-function coupling. **Left:** the global Mantel correlation between cosine and Pearson matrices stays near zero across layers. **Right:** functional connectivity enrichment at the top  $k\%$  of structural edges, normalized by the global mean. The top 0.1% shows approximately 20–25-fold enrichment in both models.

234 the total coupling mass from a given upstream community. This measure follows a non-monotonic  
 235 trajectory across depth (Figure 6). In early layers, `concentration_80` values are low (1–3% for  
 236 270M layers 0–2, 2.5–3% for 1B layers 0–1), indicating that each upstream community directs  
 237 its signal toward a small set of specific downstream targets. Through the first half of the network,  
 238 `concentration_80` rises steadily, reaching peak values of 14.4% (270M, layer 9→10) and 19.0%  
 239 (1B, layer 11→12). At these peaks, each community’s coupling spreads broadly across many  
 240 downstream targets, reflecting a diffuse, broadcast-like communication mode. In the final third  
 241 of the network, `concentration_80` decreases to 5–9% (270M) and 7.5–13.5% (1B), indicating  
 242 reconvergence toward focused coupling.

243 This inverted-U trajectory is consistent with three regimes of residual stream coupling. In early layers,  
 244 weight overlap is narrow and feature-specific, suggesting focused communication pathways. In middle  
 245 layers, coupling spreads broadly across downstream communities, consistent with a broadcast-like  
 246 mode of information distribution [11]. In late layers, coupling reconverges toward fewer downstream  
 247 targets, suggesting a return to focused processing. Whether this weight-space pattern corresponds  
 248 to actual information flow during processing remains to be tested through causal interventions. The  
 249 pattern is quantitatively consistent across both model scales.

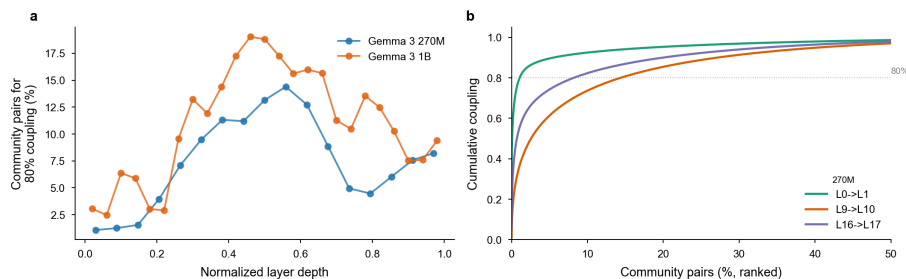


Figure 6: Cross-layer coupling trajectory. **Left:** `concentration_80` versus layer depth, showing the inverted-U shape. **Right:** cumulative coupling strength for representative layer transitions.

#### 250 4.5 Do topological structures carry interpretable semantics?

251 Whether topological structures carry genuine semantic content is a crucial test of these networks’  
 252 functional relevance. For each layer, we identify the top-20 features by degree in the density-matched  
 253 graph and examine their maximally activating corpus examples. Hub semantics follow a systematic  
 254 three-phase evolution across depth in both models. In Gemma 3 270M, the input layer is dominated  
 255 by domain-specific content words (hub degree about  $6\times$  the mean), early layers (1–3) by format and  
 256 code mega-hubs such as punctuation and structural tokens (hub degree  $34\text{--}47\times$  the mean), middle

257 layers (9–14) by mathematical operators and digits (hub degree  $11\text{--}17\times$  the mean), and the output  
 258 layer by morphological roots (hub degree  $23\times$  the mean). The 1B model exhibits a qualitatively  
 259 similar progression with shifted indices fitting its 26 layers. The mega-hub phenomenon at the format  
 260 and code layers is the network manifestation of universal-connector tokens that co-occur with almost  
 261 everything else in the corpus and therefore form the backbone of the early-layer co-activation graph.  
 262 Figure 7 shows representative hubs and LRG communities for 270M, with the parallel 1B figures in  
 263 Appendix D.

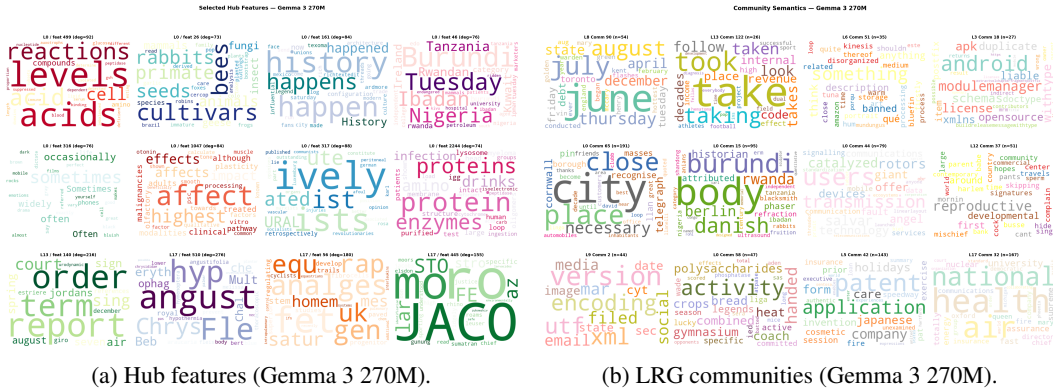


Figure 7: Semantic interpretation of topology in Gemma 3 270M. Each panel is a word cloud from the maximally activating contexts of a hub feature or representative members of a community. Selected hubs progress from domain content words at the input layer through format and code mega-hubs and digits to morphological roots at the output. Communities range from geography and technology at layer 0 through temporal and file-format clusters in middle layers to institutional terms in late layers.

264 LRG communities provide a complementary view at the module level. Features within a commu-  
 265 nity respond to the same semantic domain, syntactic category, or format type. Selected examples  
 266 include geography clusters (layer 0), temporal terms (layer 8), biology (layer 12), and morphological  
 267 paradigms (layer 13) in 270M, with comparable patterns in 1B (Appendix D). The community hierar-  
 268 chy across layers follows established models of transformer processing stages. Layer 0 communities  
 269 cluster by semantic topic and early layers by syntactic categories, while middle layers shift toward  
 270 grammatical function and late layers toward word-formation and morphological structure. This  
 271 semantic alignment between topological structure and known linguistic processing stages [21] is  
 272 consistent with the interpretation that the network properties reflect functional organization, though  
 273 high-degree features may be interpretable simply because they are high-frequency.

## 274 5 Discussion

275 Our analyses paint a unified multi-scale picture of LLM feature network organization. Within  
 276 modules at path distances below the crossover length  $l_B^* \approx 6$ , networks exhibit fractal-like scaling  
 277 with  $d_B \approx 1.35$ . Across modules at larger scales, small-world shortcuts break the fractal scaling  
 278 and produce  $\omega \approx 0$  topology. Each layer possesses a single dominant structural scale, as revealed  
 279 by the single peak in spectral heat capacity, and this scale becomes sharper with increasing depth.  
 280 Between layers, information flow through the residual stream follows a focused-to-broadcast-to-  
 281 focused trajectory. Throughout, topological structures (hubs and communities) carry interpretable  
 282 semantics aligned with known processing hierarchies.

283 This organizational profile parallels biological neural networks in several respects. Both systems  
 284 exhibit small-world topology, increasing modularity along processing hierarchies, fractal geometry  
 285 within functional modules, disproportionate vulnerability to targeted hub removal, and structure-  
 286 function dissociation coupled with strong local correspondence. The comparison is quantitative (Ta-  
 287 ble 1). The main quantitative difference is the lower fractal dimension in LLM networks ( $d_B \approx 1.35$   
 288 vs  $\approx 2.0$  in brain networks), suggesting that feature communities have sparser internal connectivity  
 289 than brain functional modules. This difference may reflect the lower-dimensional geometry of residual  
 290 stream representations compared to the three-dimensional embedding of biological cortex, or it may  
 291 arise from differences in how co-activation is defined in the two domains.

Table 1: Comparison of organizational properties between biological brain networks and LLM feature networks.

Property	Brain networks	LLM feature networks
Small-world topology	$\omega \approx 0$ [4]	$\omega \approx 0$ (this work)
Increasing modularity	primary $\rightarrow$ association cortex	early $\rightarrow$ late layers
Fractal within modules	$d_B \approx 2.0$ [12]	$d_B \approx 1.35$ (this work)
Hub vulnerability	hub lesions cause disproportionate damage	hub removal collapses LCC
SC-FC dissociation	DTI $\neq$ fMRI but locally correlated	cosine $\neq$ Pearson but locally enriched
Single characteristic scale	$C(\tau)$ single peak	$C(\tau)$ single peak (this work)

292 The inverted-U coupling pattern carries specific implications for the “residual stream as commu-  
 293 nication bus” perspective in mechanistic interpretability [11]. If early layers perform focused,  
 294 feature-specific processing and middle layers broadcast information broadly, then interventions on  
 295 the residual stream should have qualitatively different effects depending on depth. Ablating residual  
 296 stream directions in early layers should disrupt specific downstream features, while equivalent inter-  
 297 ventions in middle layers should produce more diffuse, distributed effects. This prediction is directly  
 298 testable through activation patching experiments.

299 The quantitative consistency between 270M and 1B models strengthens these conclusions. Fractal  
 300 dimensions (1.37 vs 1.34), crossover lengths (both approximately 6), single-peak heat capacity,  
 301 inverted-U coupling trajectories, and semantic evolution patterns are all qualitatively and quantita-  
 302 tively similar across a fourfold increase in parameter count. This consistency suggests that these  
 303 organizational principles are robust within the Gemma 3 family across a fourfold increase in parameter  
 304 count. Whether they generalize to other model families remains an open empirical question.

305 Several limitations warrant caution. Co-activation networks capture statistical associations, not causal  
 306 circuits. The topological centrality of hub features does not necessarily imply causal importance for  
 307 model behavior, and feature ablation experiments would be needed to test whether degree centrality  
 308 predicts functional impact. All graph metrics remain threshold-dependent. While density-matching  
 309 and multi-threshold robustness checks mitigate this concern, the small-world property in particular  
 310 is sensitive to density, holding at  $d = 0.001$  but weakening at  $d = 0.005$ . Skip-transcoders explain  
 311 approximately 85–95% of MLP variance, and the unexplained variance may carry additional relevant  
 312 structure. The greedy box-covering algorithm provides an upper bound on  $N_B$ , potentially biasing  
 313 fractal dimension estimates. Semantic interpretation through max-activating examples from 4,096  
 314 Pile samples may miss rare activation patterns, as evidenced by the 9% empty rate for 270M layer 7  
 315 hubs. Finally, our analysis covers two scales from one model family (Gemma 3). Establishing  
 316 whether these organizational principles are universal properties of language models or specific to the  
 317 Gemma architecture requires investigation across diverse model families, training procedures, and  
 318 scales.

## 319 6 Conclusion

320 We have presented, to our knowledge, the first systematic network-science analysis of an LLM’s  
 321 interpretable feature space at corpus scale, constructing and analyzing layerwise co-activation and  
 322 structural networks from Gemma 3 transcoder features. The networks exhibit analogues of all  
 323 four classical cortical signatures and additionally show a dual-regime fractal scaling, a single-  
 324 peaked spectral heat capacity that sharpens with depth, and a non-monotonic cross-layer coupling  
 325 trajectory. Hub features and coarse-grained communities carry interpretable semantics organized  
 326 along a processing hierarchy that matches known transformer computation stages. The similarity  
 327 to cortical organizational principles, consistent across two model scales, indicates that trained  
 328 transformer representations share several statistical network properties with cortical networks, though  
 329 whether this reflects convergent optimization pressures or properties common to heterogeneous  
 330 modular networks remains to be established. This work provides a foundation for a network-theoretic  
 331 approach to LLM interpretability and demonstrates that the toolbox of network neuroscience can be  
 332 productively applied to artificial information processing systems.

## 333 References

- 334 [1] Albert, R., Jeong, H. and Barabási, A.L. Error and attack tolerance of complex networks.  
335 *Nature*, 406(6794):378–382, July 2000. ISSN 0028-0836, 1476-4687. doi: 10.1038/35019019.
- 336 [2] Ameisen, E., Lindsey, J., Pearce, A., Gurnee, W., Turner, N.L., Chen, B., Citro, C., Abrahams,  
337 D., Carter, S., Hosmer, B., Marcus, J., Sklar, M., Templeton, A., Bricken, T., McDougall, C.,  
338 Cunningham, H., Henighan, T., Jermyn, A., Jones, A., Persic, A., Qi, Z., Ben Thompson, T.,  
339 Zimmerman, S., Rivoire, K., Conerly, T., Olah, C. and Batson, J. Circuit tracing: Revealing  
340 computational graphs in language models. *Transformer Circuits Thread*, 2025. URL <https://transformer-circuits.pub/2025/attribution-graphs/methods.html>.  
341
- 342 [3] Balcells, D., Lerner, B., Oesterle, M., Ucar, E. and Heimersheim, S. Evolution of sae features  
343 across layers in llms, November 2024.
- 344 [4] Bassett, D.S. and Bullmore, E.T. Small-world brain networks revisited. *The Neuroscientist*, 23  
345 (5):499–516, October 2017. ISSN 1073-8584, 1089-4098. doi: 10.1177/1073858416667720.
- 346 [5] Betzel, R.F. and Bassett, D.S. Multi-scale brain networks. *NeuroImage*, 160:73–83, October  
347 2017. ISSN 10538119. doi: 10.1016/j.neuroimage.2016.11.006.
- 348 [6] Clauset, A., Shalizi, C.R. and Newman, M.E.J. Power-law distributions in empirical data.  
349 *SIAM Review*, 51(4):661–703, November 2009. ISSN 0036-1445, 1095-7200. doi: 10.1137/  
350 070710111.
- 351 [7] Colizza, V., Flammini, A., Serrano, M.A. and Vespignani, A. Detecting rich-club ordering in  
352 complex networks. *Nature Physics*, 2(2):110–115, February 2006. ISSN 1745-2473, 1745-2481.  
353 doi: 10.1038/nphys209.
- 354 [8] Cunningham, H., Ewart, A., Riggs, L., Huben, R. and Sharkey, L. Sparse autoencoders find  
355 highly interpretable features in language models, October 2023.
- 356 [9] Deng, R., Hu, X., Gilberti, M., Storks, S., Taxali, A., Angstadt, M., Sripada, C. and Chai, J.  
357 Sparse feature coactivation reveals causal semantic modules in large language models, April  
358 2026.
- 359 [10] Dunefsky, J., Chlenski, P. and Nanda, N. Transcoders find interpretable llm feature circuits,  
360 November 2024.
- 361 [11] Elhage, N., Nanda, N., Olsson, C., Henighan, T., Joseph, N., Mann, B., Askell, A., Bai, Y.,  
362 Chen, A., Conerly, T., DasSarma, N., Drain, D., Ganguli, D., Hatfield-Dodds, Z., Hernandez,  
363 D., Jones, A., Kernion, J., Lovitt, L., Ndousse, K., Amodei, D., Brown, T., Clark, J., Kaplan, J.,  
364 McCandlish, S. and Olah, C. A mathematical framework for transformer circuits. *Transformer*  
365 *Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- 366 [12] Gallos, L.K., Makse, H.A. and Sigman, M. A small world of weak ties provides optimal global  
367 integration of self-similar modules in functional brain networks. *Proceedings of the National*  
368 *Academy of Sciences*, 109(8):2825–2830, February 2012. ISSN 0027-8424, 1091-6490. doi:  
369 10.1073/pnas.1106612109.
- 370 [13] Humphries, M.D. and Gurney, K. Network ‘small-world-ness’: A quantitative method for  
371 determining canonical network equivalence. *PLoS ONE*, 3(4):e0002051, April 2008. ISSN  
372 1932-6203. doi: 10.1371/journal.pone.0002051.
- 373 [14] Kashtan, N. and Alon, U. Spontaneous evolution of modularity and network motifs. *Proceedings*  
374 *of the National Academy of Sciences*, 102(39):13773–13778, September 2005. ISSN 0027-8424,  
375 1091-6490. doi: 10.1073/pnas.0503610102.
- 376 [15] Kornblith, S., Norouzi, M., Lee, H. and Hinton, G.E. Similarity of neural network represen-  
377 tations revisited. In *Proceedings of the 36th International Conference on Machine Learning*,  
378 volume 97 of *Proceedings of Machine Learning Research*, pages 3519–3529. PMLR, 2019.
- 379 [16] Kriegeskorte, N., Mur, M. and Bandettini, P. Representational similarity analysis – connecting  
380 the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2:4, 2008. doi:  
381 10.3389/neuro.06.004.2008.

- 382 [17] Laptev, D., Balagansky, N., Aksenov, Y. and Gavrilov, D. Analyze feature flow to enhance  
383 interpretation and steering in language models, July 2025.
- 384 [18] Lieberum, T., Rajamanoharan, S., Conmy, A., Smith, L., Sonnerat, N., Varma, V., Kramár, J.,  
385 Dragan, A., Shah, R. and Nanda, N. Gemma scope: Open sparse autoencoders everywhere all  
386 at once on gemma 2, August 2024.
- 387 [19] Liégeois, R., Santos, A., Matta, V., Van De Ville, D. and Sayed, A.H. Revisiting correlation-  
388 based functional connectivity and its relationship with structural connectivity. *Network Neuro-*  
389 *science*, 4(4):1235–1251, January 2020. ISSN 2472-1751. doi: 10.1162/netn\_a\_00166.
- 390 [20] Lin, B. and Kriegeskorte, N. The topology and geometry of neural representations. *Proceedings*  
391 *of the National Academy of Sciences*, 121(42):e2317881121, October 2024. ISSN 0027-8424,  
392 1091-6490. doi: 10.1073/pnas.2317881121.
- 393 [21] Lindsey, J., Gurnee, W., Ameisen, E., Chen, B., Pearce, A., Turner, N.L., Citro, C., Abra-  
394 hams, D., Carter, S., Hosmer, B., Marcus, J., Sklar, M., Templeton, A., Bricken, T., Mc-  
395 Dougall, C., Cunningham, H., Henighan, T., Jermyn, A., Jones, A., Persic, A., Qi, Z.,  
396 Thompson, T.B., Zimmerman, S., Rivoire, K., Conerly, T., Olah, C. and Batson, J. On  
397 the biology of a large language model. *Transformer Circuits Thread*, 2025. URL <https://transformer-circuits.pub/2025/attribution-graphs/biology.html>.  
398
- 399 [22] Liu, Y., Liu, Z., Wu, Z., Ning, J., Sun, H., Xia, S., Yang, Y., Gao, X., Qiang, N., Ge, B., Liu, T.,  
400 Han, J. and Hu, X. Brain-inspired exploration of functional networks and key neurons in large  
401 language models, January 2026.
- 402 [23] Marks, S., Rager, C., Michaud, E.J., Belinkov, Y., Bau, D. and Mueller, A. Sparse feature  
403 circuits: Discovering and editing interpretable causal graphs in language models, March 2025.
- 404 [24] McDougall, C., Conmy, A., Kramár, J., Lieberum, T., Rajamanoharan, S. and Nanda, N. Gemma  
405 Scope 2 - technical paper. Technical report, Google, 2025.
- 406 [25] Mocanu, D.C., Mocanu, E., Stone, P., Nguyen, P.H., Gibescu, M. and Liotta, A. Scalable  
407 training of artificial neural networks with adaptive sparse connectivity inspired by network  
408 science. *Nature Communications*, 9(1):2383, June 2018. ISSN 2041-1723. doi: 10.1038/  
409 s41467-018-04316-3.
- 410 [26] Newman, M.E.J. Assortative mixing in networks. *Physical Review Letters*, 89(20):208701,  
411 October 2002. ISSN 0031-9007, 1079-7114. doi: 10.1103/PhysRevLett.89.208701.
- 412 [27] Newman, M.E.J. Modularity and community structure in networks. *Proceedings of the National*  
413 *Academy of Sciences*, 103(23):8577–8582, June 2006. ISSN 0027-8424, 1091-6490. doi:  
414 10.1073/pnas.0601602103.
- 415 [28] Newman, M.E.J. *Networks*. Oxford University Press, Oxford, United Kingdom ; New York,  
416 NY, United States of America, second edition edition, 2018. ISBN 978-0-19-880509-0.
- 417 [29] Paulo, G., Shabalín, S. and Belrose, N. Transcoders beat sparse autoencoders for interpretability,  
418 February 2025.
- 419 [30] Penedo, G., Kydlíček, H., allal, L.B., Lozhkov, A., Mitchell, M., Raffel, C., Werra, L.V. and  
420 Wolf, T. The fineweb datasets: Decanting the web for the finest text data at scale, October 2024.
- 421 [31] Poggialini, A., Villegas, P., Muñoz, M.A. and Gabrielli, A. Networks with many structural  
422 scales: A renormalization group perspective, December 2024.
- 423 [32] Schrimpf, M., Blank, I.A., Tuckute, G., Kauf, C., Hosseini, E.A., Kanwisher, N., Tenenbaum,  
424 J.B. and Fedorenko, E. The neural architecture of language: Integrative modeling converges on  
425 predictive processing. *Proceedings of the National Academy of Sciences*, 118(45):e2105646118,  
426 2021. doi: 10.1073/pnas.2105646118.
- 427 [33] Song, C., Havlin, S. and Makse, H.A. Self-similarity of complex networks. *Nature*, 433(7024):  
428 392–395, January 2005. ISSN 0028-0836, 1476-4687. doi: 10.1038/nature03248.

- 429 [34] Song, C., Gallos, L.K., Havlin, S. and Makse, H.A. How to calculate the fractal dimension of  
430 a complex network: The box covering algorithm. *Journal of Statistical Mechanics: Theory*  
431 *and Experiment*, 2007(03):P03006–P03006, March 2007. ISSN 1742-5468. doi: 10.1088/  
432 1742-5468/2007/03/P03006.
- 433 [35] Sporns, O. Network attributes for segregation and integration in the human brain. *Current*  
434 *Opinion in Neurobiology*, 23(2):162–171, April 2013. ISSN 09594388. doi: 10.1016/j.conb.  
435 2012.11.015.
- 436 [36] Team, G., Kamath, A., Ferret, J., Pathak, S., Vieillard, N., Merhej, R., Perrin, S., Matejovicova,  
437 T., Ramé, A., Rivière, M., Rouillard, L., Mesnard, T., Cideron, G., bastien Grill, J., Ramos,  
438 S., Yvinec, E., Casbon, M., Pot, E., Penchev, I., Liu, G., Visin, F., Kenealy, K., Beyer, L.,  
439 Zhai, X., Tsitsulin, A., Busa-Fekete, R., Feng, A., Sachdeva, N., Coleman, B., Gao, Y.,  
440 Mustafa, B., Barr, I., Parisotto, E., Tian, D., Eyal, M., Cherry, C., Peter, J.T., Sinopalnikov,  
441 D., Bhupatiraju, S., Agarwal, R., Kazemi, M., Malkin, D., Kumar, R., Vilar, D., Brusilovsky,  
442 I., Luo, J., Steiner, A., Friesen, A., Sharma, A. et al. Gemma 3 technical report, 2025. URL  
443 <https://arxiv.org/abs/2503.19786>.
- 444 [37] Telesford, Q.K., Joyce, K.E., Hayasaka, S., Burdette, J.H. and Laurienti, P.J. The ubiquity of  
445 small-world networks. *Brain Connectivity*, 1(5):367–375, December 2011. ISSN 2158-0014,  
446 2158-0022. doi: 10.1089/brain.2011.0038.
- 447 [38] Templeton, A., Conerly, T., Marcus, J., Lindsey, J., Bricken, T., Chen, B., Pearce, A., Citro,  
448 C., Ameisen, E., Jones, A., Cunningham, H., Turner, N.L., McDougall, C., MacDiarmid,  
449 M., Freeman, C.D., Summers, T.R., Rees, E., Batson, J., Jermyn, A., Carter, S., Olah, C. and  
450 Henighan, T. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet.  
451 *Transformer Circuits Thread*, 2024. URL [https://transformer-circuits.pub/2024/  
452 scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).
- 453 [39] Traag, V.A., Waltman, L. and Van Eck, N.J. From louvain to leiden: Guaranteeing well-  
454 connected communities. *Scientific Reports*, 9(1):5233, March 2019. ISSN 2045-2322. doi:  
455 10.1038/s41598-019-41695-z.
- 456 [40] Villegas, P., Gili, T., Caldarelli, G. and Gabrielli, A. Laplacian renormalization group for  
457 heterogeneous networks. *Nature Physics*, 19(3):445–450, March 2023. ISSN 1745-2473,  
458 1745-2481. doi: 10.1038/s41567-022-01866-8.
- 459 [41] Watts, D.J. and Strogatz, S.H. Collective dynamics of ‘small-world’ networks. *Nature*, 393  
460 (6684):440–442, June 1998. ISSN 0028-0836, 1476-4687. doi: 10.1038/30918.
- 461 [42] Yamins, D.L.K. and DiCarlo, J.J. Using goal-driven deep learning models to understand sensory  
462 cortex. *Nature Neuroscience*, 19(3):356–365, 2016. doi: 10.1038/nn.4244.
- 463 [43] Zheng, Y., Yuan, Y., Zhuo, Y., Li, Y., Kreiman, G., Poggio, T. and Santi, P. Probing neural  
464 topology of large language models, January 2026.

465 **A Preliminaries**

466 A decoder-only transformer is a stack of identical blocks, each composed of a multi-head self-attention  
 467 sublayer and a position-wise MLP sublayer that read from and write to a shared residual stream  
 468 of dimension  $d_{\text{model}}$  [11]. The MLP sublayers carry most of the parameter budget and most of the  
 469 compositional computation, so mechanistic interpretability work has typically located the units we  
 470 informally call “features” inside or around them.

471 A sparse autoencoder (SAE) approximates an activation  $x \in \mathbb{R}^{d_{\text{model}}}$  as a sparse combination of  
 472 learned dictionary atoms through an encoder  $z = \phi(W_{\text{enc}}x + b_{\text{enc}}) \in \mathbb{R}^{d_{\text{sae}}}$  and a linear decoder  
 473  $\hat{x} = W_{\text{dec}}z + b_{\text{dec}}$  with  $d_{\text{sae}} \gg d_{\text{model}}$  [8, 38]. Each column of  $W_{\text{dec}}$  is a feature direction in the  
 474 residual stream and each entry of  $z$  is its activation on the current token. A transcoder generalises this  
 475 to the input-output map of an MLP block [10, 29], and a skip transcoder [24] adds a learned linear  
 476 shortcut,

$$\widehat{\text{MLP}}(x) = W_{\text{dec}} z(x) + b_{\text{dec}} + W_{\text{skip}} x, \quad (1)$$

477 where  $W_{\text{skip}}x$  absorbs the purely linear part and  $z(x)$  carries the feature-like residual. The transcoders  
 478 we use are 16,384-dimensional with JumpReLU activations and an L0 of about 50 active features per  
 479 token, yielding a sparsity of approximately 0.3%.

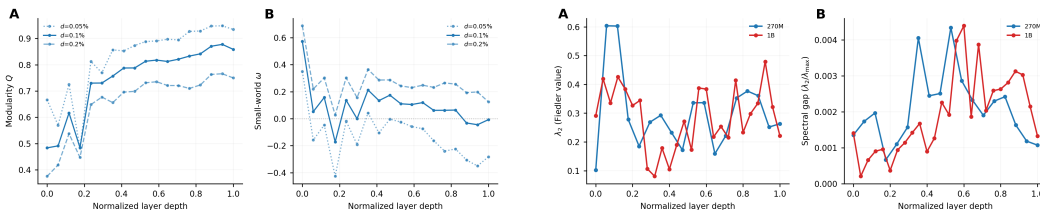
480 Treating each (layer, feature) pair as a node and weighting edges by co-activation across a fixed corpus,  
 481 we summarise the resulting graph using the standard network-science quantities. These include the  
 482 clustering coefficient  $C$  and average shortest path length  $L$  combined into the Telesford small-world  
 483 index  $\omega = L_{\text{rand}}/L - C/C_{\text{lattice}}$  [37], the Leiden modularity  $Q$  [39], the degree assortativity  $r$  [26],  
 484 percolation under random and targeted node removal [1], and Mantel correlations between weight-  
 485 based and activation-based matrices following the SC-FC methodology of Liégeois et al. [19]. For  
 486 multi-scale structure we add three renormalisation group tools, namely box covering for the fractal  
 487 dimension  $d_B$  [34, 12], the Laplacian renormalisation group of Villegas et al. [40], and the spectral  
 488 heat capacity  $C(\tau)$  of Poggialini et al. [31]. This toolbox is standard in network neuroscience, which  
 489 is what makes the cortex-LLM comparison meaningful.

490 **B Threshold-independent and Robustness Checks**

491 This appendix collects additional figures and tables that support claims made in the main text. Table 2  
 492 reports per-model summaries of the canonical density-matched graph statistics. Figure 8a shows the  
 493 modularity and small-world index at three density levels for Gemma 3 270M, Figure 8b shows the  
 494 threshold-free Laplacian spectrum, and Figures 9a and 9b show degree distribution properties and the  
 495 Configuration Model significance heatmap.

Table 2: Layerwise summary of density-matched ( $d = 0.001$ ) topology metrics across all layers.

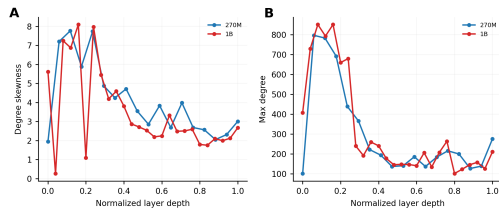
Model	$C$	$L$	$\omega$	$Q$ (Leiden)	$r$	deg. skew
Gemma 3 270M	0.25 to 0.46	4.0 to 7.7	-0.17 to +0.57	0.48 to 0.88	0.18 to 0.78	1.9 to 7.8
Gemma 3 1B	0.30 to 0.83	1.9 to 7.6	-0.43 to +0.32	0.09 to 0.93	-0.43 to 0.82	0.3 to 8.1



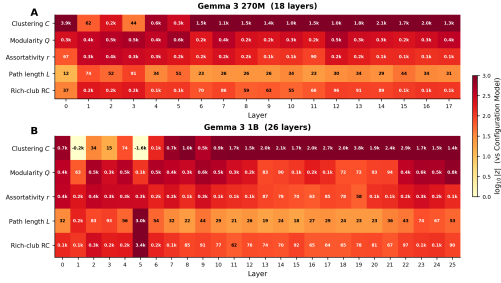
(a) Density robustness ( $d = 0.0005, 0.001, 0.002$ ).

(b) Threshold-free Laplacian spectrum.

Figure 8: Robustness checks. (a) Modularity  $Q$  and small-world  $\omega$  at three density levels for Gemma 3 270M, showing depth trends consistent across densities. (b) Fiedler value  $\lambda_2$  and normalised spectral gap on unthresholded Pearson matrices, which display the same depth dependence as the thresholded metrics without any threshold choice.



(a) Degree skewness and maximum degree.

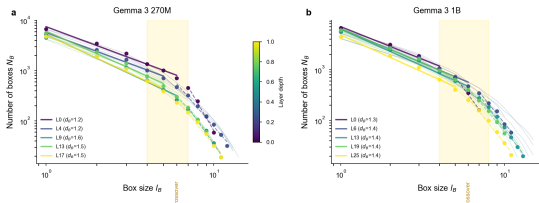


(b) Configuration Model z-score heatmap.

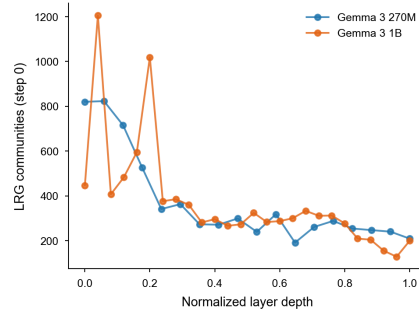
Figure 9: Degree distribution and Configuration Model significance. (a) Degree skewness and maximum degree as a function of layer depth, with the two anomalous layers in 1B (layer 1 and layer 5) visible as outliers. (b) Configuration Model z-score heatmap for clustering and modularity across all layers and density levels, with almost every cell in the high-significance range.

## 496 C Renormalisation Group Supplementary Plots

497 Figure 10 collects the per-layer log-log box-covering curves and the Laplacian renormalisation group  
 498 supernode count.



(a) Per-layer log-log box-covering curves with piecewise fits.



(b) LRG supernode count vs depth.

Figure 10: Renormalisation group supplementary plots. (a) Box-covering scaling curves for representative layers, with piecewise linear fits and the crossover band shaded around  $l_B^* \approx 6$ . (b) Laplacian renormalisation group supernode count as a function of layer depth. The monotonic decrease shows progressively more compressible structure in deeper layers.

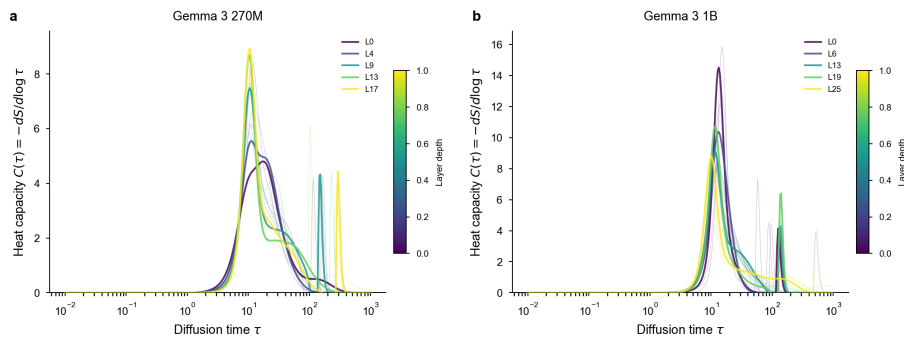


Figure 11: Spectral heat capacity  $C(\tau)$  for each layer of Gemma 3 270M and 1B. All curves show a single peak (no plateau), ruling out scale invariance. Peaks shift leftward and grow taller with depth.

499 **D Semantic Interpretation for Gemma 3 1B**

500 Figure 12 shows the parallel hub and community word-cloud panels for Gemma 3 1B. The semantic  
 501 patterns are qualitatively the same as in 270M (Figure 7) but with shifted layer indices that fit the 26  
 502 layers of 1B.

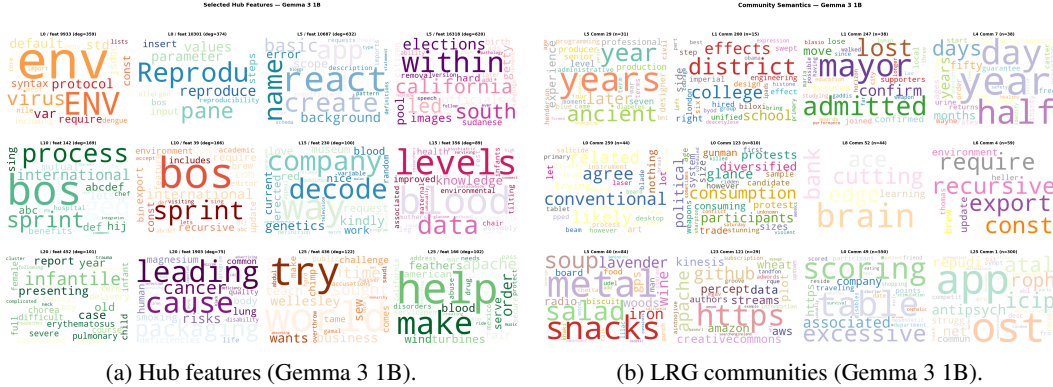


Figure 12: Semantic interpretation in Gemma 3 1B. The depth-dependent semantic pattern matches 270M, with HTML and multilingual tokens taking the early role and Latin and morphological roots taking the late role.

503 **E Extended Methods**

504 **Network construction details.** A forward hook on every MLP module captures the MLP input on  
 505 each forward pass and feeds it through the corresponding skip transcoder, producing a sparse vector  
 506 in  $\mathbb{R}^{16,384}$ . With about 50 active features per token, only  $\binom{50}{2} \approx 1225$  feature pairs are updated  
 507 per token, compared to the full  $1.34 \times 10^8$  pairs in a 16K square matrix. We accumulate sufficient  
 508 statistics for Pearson, Jaccard, cosine, and NPMI correlations in a single streaming pass without  
 509 storing raw activations, and the resulting sufficient statistics are computed exactly (without streaming  
 510 approximation). The output per layer is a 16K by 16K co-activation matrix together with a vector  
 511 of per-feature activation counts used as a quality filter. Features with activation count below ten are  
 512 excluded before any downstream analysis.

513 For the structural networks, taking absolute values of the virtual weight  $M_{ij}^{(\ell \rightarrow \ell+1)}$  is a deliberate  
 514 choice that aligns the structural network with the unsigned co-activation network and keeps edge  
 515 weights compatible with modularity, box covering, and Mantel statistics, all of which require non-  
 516 negative weights. The signed driving relation between features, which would distinguish inhibition  
 517 from excitation, is left to follow-up work. The 270M model has  $d_{\text{model}} = 640$ , which makes its  
 518 within-layer cosine similarity matrix too low in contrast for meaningful thresholding, so for 270M we  
 519 report only cross-layer virtual weights as structural connectivity.

520 **Density-matched thresholding.** Graphs are obtained from the dense matrices by retaining the top  
 521 edges by absolute correlation until the graph reaches a target density. Density matching is preferred  
 522 over a fixed threshold because the scale of the correlation matrix varies systematically with depth.  
 523 The canonical density is  $d = 0.001$ , with  $d = 0.0005$  and  $d = 0.002$  as robustness checks. A  
 524 fixed-threshold control at  $\theta = 0.15$  is also reported in the robustness analysis. All graph metrics are  
 525 computed on the largest connected component using igraph.

526 **Topology metrics.** The full set of computed metrics includes the clustering coefficient  $C$ , the  
 527 average shortest path length  $L$ , the Telesford small-world index  $\omega = L_{\text{rand}}/L - C/C_{\text{lattice}}$  [37], the  
 528 Humphries-Gurney index  $\sigma$  [13], Leiden modularity  $Q$  at resolutions  $\gamma \in \{0.5, 1.0, 1.5, 2.0\}$  [39],  
 529 degree skewness, the Clauset-Shalizi-Newman power-law and lognormal fits to the degree distribution  
 530 [6], degree assortativity  $r$  [26], and the rich-club coefficient at the median degree [7].

531 **Null models.** Every metric is tested against two null models using fifty realisations per layer.  
532 The Erdős-Rényi null preserves the number of nodes and edges. The Configuration Model null  
533 additionally preserves the degree sequence [28], providing a more demanding baseline because any  
534 property that survives this comparison cannot be explained by the degree distribution alone.

535 **Percolation protocol.** Percolation removes 0 to 50 percent of nodes under three strategies and  
536 tracks the fraction of the largest connected component remaining. Random removal averages  
537 over twenty independent trials. Degree-targeted removal recalculates degrees after each step, and  
538 betweenness-targeted removal recalculates betweenness centrality after each step. A Configuration  
539 Model percolation control is run in parallel to separate genuine topological vulnerability from effects  
540 that follow from the degree distribution alone.

541 **Structure-function coupling.** Following Liégeois et al. [19], we report three complementary  
542 measures. The Mantel correlation between the structural and functional matrices uses node-label  
543 permutation to generate the null distribution, which tests whether structurally connected feature pairs  
544 have higher co-activation than expected. The conditional analysis computes the mean functional  
545 connectivity inside the top  $k$  percent of structural edges for  $k \in \{0.1, 0.5, 1, 5, 10\}$ , normalised by the  
546 global mean, to test whether alignment concentrates at the strongest structural ties. A threshold-free  
547 check uses the Fiedler value  $\lambda_2$  and the normalised spectral gap of the unthresholded Pearson matrix,  
548 which display the same depth dependence as the thresholded metrics without any threshold choice.

549 **Renormalization group details.** For box covering, we run the greedy algorithm of Song et al.  
550 [34] on each layer’s largest connected component, sweeping the box diameter  $l_B$  from 2 to the  
551 graph diameter and fitting the scaling law  $N_B(l_B) \sim l_B^{-d_B}$ . To test for dual-regime scaling, we  
552 additionally fit a piecewise linear model in log-log space with a breakpoint  $l_B^*$  chosen by minimising  
553 the total residual sum of squares. The piecewise fit is accepted only when  $\Delta\text{AIC} > 10$  relative to  
554 the single-regime fit, and the local exponent is bootstrapped to obtain a 95% confidence interval.  
555 The Laplacian renormalization group of Villegas et al. [40] builds the normalised Laplacian  $L_{\text{norm}} =$   
556  $I - D^{-1/2}AD^{-1/2}$  and the diffusion kernel  $K(\tau) = \exp(-\tau L_{\text{norm}})$ , then merges nodes that share  
557 diffusive modes by agglomerative clustering on the rows of  $K(\tau)$ . The compression ratio, defined as  
558 original nodes divided by supernodes, measures structural redundancy at scale  $\tau$ .

559 **Semantic interpretation pipeline.** For each layer we select the top twenty hub features by degree  
560 in the canonical  $d = 0.001$  graph and feed 4,096 samples from Pile-uncopyrighted through the  
561 hook-based pipeline, recording the top ten (token, context, activation) triples per hub. For Laplacian  
562 renormalization group communities of size at least five, we take the top fifty communities per  
563 layer and pick five representative features by within-community weighted degree, running the same  
564 activation-recording pipeline for each.

565 **Reproducibility.** The streaming co-activation tracker is in  
566 `src/network_construction/activation_based.py`. The weight-based and cross-  
567 layer network builders are in `src/network_construction/weight_based.py`. Topol-  
568 ogy analysis scripts (`compute_metrics.py`, `percolation.py`, `spectral_analysis.py`,  
569 `structure_function_coupling.py`, `null_model_significance.py`) and renormalization  
570 group scripts (`box_covering.py`, `lrg.py`, `heat_capacity.py`) reproduce every number reported  
571 in this paper. Hub and community semantic labels are produced by `hub_interpretation.py` and  
572 `community_interpretation.py`. All experiments run in a single environment on an NVIDIA  
573 RTX 5090 with CUDA 12.8 and a recent PyTorch build.