

1 Understanding Oversmoothing in Deep Graph Neural Networks: Insights from
2 Sparsity, Weights, and Lottery Tickets

3 Ana Carolina Ramos Cunha

4 February 22, 2026

5 **1 Introduction**

6 Graph Neural Networks (GNNs) have become a fundamental tool for learning on graph-structured
7 data, achieving strong performance in node classification, link prediction, and graph-level tasks.
8 Architectures such as Graph Convolutional Networks (GCNs), Graph Attention Networks (GATs),
9 and Graph Isomorphism Networks (GINs) rely on iterative message passing to aggregate informa-
10 tion from local neighborhoods. While effective in shallow settings, deep GNNs suffer from a critical
11 limitation known as *oversmoothing*.

12 Oversmoothing refers to the phenomenon in which node representations become increasingly indis-
13 tinguishable as the number of layers grows. Repeated neighborhood aggregation causes embeddings
14 to converge toward a low-dimensional subspace dominated by low-frequency graph components,
15 ultimately leading to representation collapse and severe performance degradation. As a result,
16 standard GNN architectures typically fail beyond 5–8 layers.

17 Recent work has suggested that oversmoothing may not be an inherent structural limitation of
18 GNNs, but rather a consequence of weight optimization dynamics. In particular, the Untrained
19 Graph Neural Networks Tickets (UGTs) framework demonstrates that sparse subnetworks with
20 randomly initialized and frozen weights can outperform fully trained dense networks in deep
21 regimes. This challenges the conventional assumption that training dense weights is necessary
22 for strong performance.

23 In this report, we reproduce the main experimental findings of the UGTs paper and investigate
24 whether sparsity, initialization, or training dynamics are the primary drivers of oversmoothing.
25 We then extend the analysis by studying additional sparsification strategies and weight reparam-
26 eterization techniques to better understand the mechanisms that enable deep GNNs to remain
27 stable.

28 **Roadmap.** Section 2 situates our work within the broader literature on oversmoothing and
29 lottery tickets in GNNs. Section 3 details our reproduction of the original UGT results. Section 4
30 presents our extensions, including Unified Graph Sparsification (UGS), Weight Reparameterization
31 (WeightRep), and initialization studies. We conclude in Section 6 with a synthesis of findings and
32 directions for future work.

33 2 Related Work

34 2.1 Oversmoothing in Graph Neural Networks

35 Oversmoothing has been extensively studied as a fundamental limitation of deep GNNs. Early
36 analyses linked the phenomenon to repeated multiplication by graph propagation operators, which
37 act as low-pass filters that suppress high-frequency components of node signals. Chen et al. (3)
38 introduced the Mean Average Distance (MAD) metric to quantify representation smoothness and
39 demonstrated that deep GNNs rapidly collapse node embeddings.

40 Subsequent theoretical work connected oversmoothing to spectral contraction and rank deficiency
41 of learned representations. These analyses suggest that depth amplifies smoothing effects inherent
42 in message passing, raising the question of whether architectural modifications or training strategies
43 can counteract this collapse.

44 2.2 Lottery Tickets and Sparsity in GNNs

45 The Lottery Ticket Hypothesis, originally proposed for dense neural networks, suggests that sparse
46 subnetworks exist within randomly initialized models that can match full-network performance.
47 Ramanujan et al. (2) showed that subnetworks selected without weight training can perform com-
48 petitively when appropriate masks are learned.

49 Huang et al. (1) extended this idea to graph neural networks through Untrained GNN Tickets
50 (UGTs). In UGTs, weights remain randomly initialized and frozen, while binary masks over
51 edges are learned via task loss. Surprisingly, such untrained sparse networks can outperform
52 fully trained dense GNNs at greater depths, suggesting that training dynamics may exacerbate
53 oversmoothing.

54 Chen et al. (4) proposed Unified Graph Sparsification (UGS), which jointly sparsifies model pa-
55 rameters and graph structure through gradient-based training. Unlike UGTs, UGS retains weight
56 optimization, providing a bridge between dense training and untrained sparsity.

57 More recently, Zhuo et al. (5) introduced Weight Reparameterization (WeightRep), construct-
58 ing input-dependent weights during training to preserve representation diversity. This approach
59 demonstrates that oversmoothing can be mitigated even under training, provided weights evolve
60 in a structurally appropriate manner.

61 2.3 Positioning Our Work

62 Our study builds on these lines of research by systematically comparing: Trained dense GNNs,
63 Untrained sparse GNNs (UGTs), Edge-Popup sparsification, Unified Graph Sparsification (UGS),
64 Weight reparameterization and initialization strategies.

65 By reproducing the core UGT findings and extending the analysis to additional sparsification and
66 weight-control mechanisms, we aim to disentangle the roles of sparsity, initialization, and training
67 in the emergence of oversmoothing.

68 3 Reproduction

69 We reproduce the central experimental results of Huang et al. (1), focusing on three key analyses:
70 1. Accuracy as a function of network depth, 2. Mean Average Distance (MAD) as a function of
71 depth and training epoch, 3. Accuracy as a function of sparsity level.

72 Our reproduction covers three architectures (GCN, GAT, GIN) and three citation datasets (Cora,
73 Citeseer, Pubmed), following the original evaluation protocol.

74 3.1 Formal Setting

75 Let $G = (V, E)$ denote a graph with adjacency matrix $A \in \mathbb{R}^{n \times n}$ and node feature matrix $X \in$
76 $\mathbb{R}^{n \times d}$. A L -layer GNN computes node representations iteratively as

$$77 \quad H^{(0)} = X, \quad (1)$$

$$78 \quad H^{(\ell+1)} = \sigma \left(\mathcal{P}(A)H^{(\ell)}W^{(\ell)} \right), \quad (2)$$

79 where: $\mathcal{P}(A)$ is a graph propagation operator (e.g., normalized adjacency for GCN), $W^{(\ell)}$ are
80 learnable weights, σ is a nonlinear activation function.

81 The model is trained by minimizing cross-entropy loss:

$$82 \quad \mathcal{L}(\theta) = \sum_{i \in \mathcal{V}_{\text{train}}} \text{CE}(f_{\theta}(X, A)_i, y_i), \quad (3)$$

83 where $\theta = \{W^{(\ell)}\}_{\ell=1}^L$.

84 Here, $f_{\theta} : \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times C}$ denotes the full L -layer GNN followed by a linear classification
85 layer, mapping node features and the adjacency matrix to class logits for each node, where C is
86 the number of classes.

87 3.2 Untrained GNN Tickets (UGTs)

88 In UGTs, weights are randomly initialized and *never updated*. Instead, each weight matrix is paired
89 with a trainable score matrix $S^{(\ell)}$.

90 A binary mask $m^{(\ell)}$ is constructed via global thresholding:

$$91 \quad m_{ij}^{(\ell)} = \mathbf{1} \left[S_{ij}^{(\ell)} \geq \tau \right], \quad (4)$$

92 where τ enforces a desired sparsity level.

93 The threshold τ is chosen such that a fixed global sparsity level $s \in [0, 1]$ is enforced across all
94 layers, meaning that only the top $(1 - s)$ fraction of scores are retained.

95 The effective weights become:

$$96 \quad \widetilde{W}^{(\ell)} = W^{(\ell)} \odot m^{(\ell)}. \quad (5)$$

97 Only the scores are updated:

$$98 \quad S \leftarrow S - \eta \nabla_S \mathcal{L}(f(X, A; W \odot m), y), \quad (6)$$

99 using a straight-through estimator for the binary mask.

100 We evaluate sparsity levels of 0.1 and 0.5, and additionally sweep sparsity from 0 to 0.99 for
 101 robustness analysis.

102 3.3 Edge-Popup

103 Edge-Popup also freezes weights and learns scores, but differs in three aspects: layer-wise top-
 104 k thresholding instead of global thresholding, one-shot sparsification instead of gradual pruning,
 105 mask construction using $|S|$ instead of raw scores.

106 We additionally implement a variation using raw scores (without absolute value), observing im-
 107 proved performance at lower sparsities.

108 3.4 Mean Average Distance (MAD)

109 To quantify representation smoothness across layers, we adopt the Mean Average Distance (MAD)
 110 metric defined in Chen et al. (AAAI 2020).

111 Let $H^{(\ell)} \in \mathbb{R}^{n \times d}$ denote the matrix of node embeddings at layer ℓ , where $H_i^{(\ell)}$ is the embedding
 112 of node i .

113 First, define the cosine distance between two node embeddings:

$$114 \quad D_{ij}^{(\ell)} = 1 - \frac{H_i^{(\ell)} \cdot H_j^{(\ell)}}{\|H_i^{(\ell)}\|_2 \|H_j^{(\ell)}\|_2}. \quad (7)$$

115 Let $M_{\text{tgt}} \in \{0, 1\}^{n \times n}$ be a target mask matrix (e.g., excluding self-pairs or restricting to specific
 116 node sets). The masked distance matrix is:

$$117 \quad \tilde{D}^{(\ell)} = M_{\text{tgt}} \odot D^{(\ell)}. \quad (8)$$

118 For each node i , define the row-wise mean over non-zero entries:

$$119 \quad \text{MAD}_i^{(\ell)} = \frac{\sum_{j=1}^n \tilde{D}_{ij}^{(\ell)}}{\sum_{j=1}^n M_{\text{tgt},ij}}. \quad (9)$$

120 The overall Mean Average Distance at layer ℓ is:

$$121 \quad \text{MAD}^{(\ell)} = \frac{1}{n} \sum_{i=1}^n \text{MAD}_i^{(\ell)}. \quad (10)$$

122 Lower MAD indicates stronger representation collapse.

123 3.5 Experimental Protocol

124 We follow the original setup: depth varied from 2 to 20 layers, MAD computed for a 32-layer
 125 GCN on Cora, accuracy averaged across multiple runs, Kaiming uniform initialization for frozen
 126 weights, identical optimizer and learning rates as the original implementation.

127 **3.6 Reproduced Results**

128 **3.6.1 Accuracy vs Depth**

129 Trained dense GNNs exhibit severe degradation beyond 5–8 layers, with accuracy dropping sharply
 130 at greater depths. In contrast, UGTs maintains stable performance up to 20 layers across archi-
 131 tectures and datasets. Sparsity 0.5 often yields the strongest deep performance.

132 Edge-Popup improves over dense training but consistently underperforms UGTs in deep regimes.

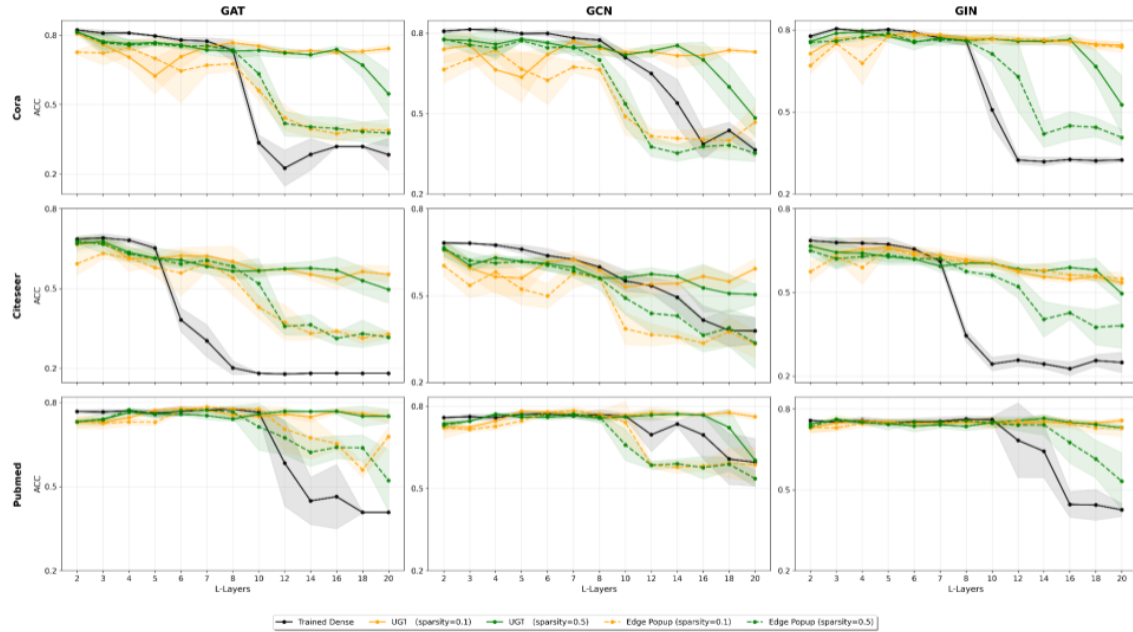


Figure 1: Accuracy versus number of layers for GCN, GAT, and GIN across datasets. Reproduced from our implementation (corresponding to Figure 2 in (1)).

133 **3.6.2 MAD Analysis**

134 Dense GCNs rapidly exhibit representation collapse, with MAD approaching zero within the first
 135 few layers. UGTs maintains significantly higher MAD across depth and training epochs, indicating
 136 preserved feature diversity.

137 These findings confirm that oversmoothing is strongly mitigated when weights remain untrained
 138 and sparsity is imposed.

139 **3.6.3 Accuracy vs Sparsity**

140 UGTs matches or slightly underperforms trained dense networks across sparsity levels from 0.1 to
 141 0.9. Performance degrades gradually at extreme sparsity (0.95–0.99).

142 Edge-Popup exhibits earlier degradation, particularly on Citeseer.

143 Notably, using raw scores instead of $|S|$ improves Edge-Popup at lower sparsities, suggesting that
 144 signed score competition enhances expressivity.

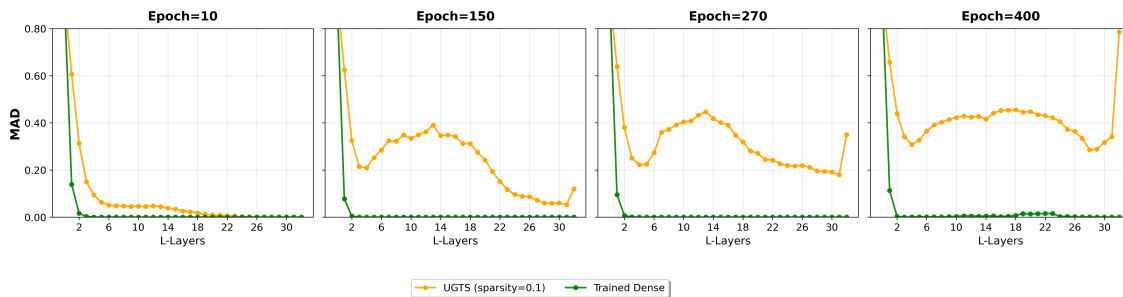


Figure 2: Mean Average Distance (MAD) across layers for a 32-layer GCN on Cora, evaluated at different training epochs. Dense training exhibits rapid representation collapse ($MAD \rightarrow 0$), while UGT maintains higher pairwise embedding distances across depth. Reproduction of Figure 4 from (1).

145 4 Extensions

146 While UGTs demonstrate that untrained sparse subnetworks can mitigate oversmoothing, several
 147 fundamental questions remain:

- 148 • Is sparsity alone sufficient to prevent oversmoothing under training?
- 149 • Can trained sparse models match the depth robustness of UGTs?
- 150 • Is oversmoothing primarily a structural issue or a learning dynamics issue?
- 151 • Can weight reparameterization or initialization strategies resolve collapse while retaining
 152 training?

153 To address these questions, we investigate three extensions: (1) Unified Graph Sparsification
 154 (UGS), (2) Weight Reparameterization (WeightRep) and Weight Initialization (WeightInit), (3)
 155 The interaction between UGT and structured initialization.

156 4.1 Unified Graph Sparsification (UGS)

157 4.1.1 Motivation

158 UGTs suggests that freezing random weights avoids representation collapse. However, this raises a
 159 key question: is the benefit due to sparsity itself, or due to the absence of weight training?

160 Unified Graph Sparsification (UGS) (4) jointly sparsifies both model weights and graph edges while
 161 retaining gradient-based weight optimization. UGS therefore separates the effects of sparsity from
 162 the absence of training.

163 4.1.2 Formal Setting

164 UGS introduces binary masks over both weights and adjacency:

$$165 \quad \widetilde{W}^{(\ell)} = W^{(\ell)} \odot m_W^{(\ell)}, \quad (11)$$

$$166 \quad \widetilde{A} = A \odot m_A, \quad (12)$$

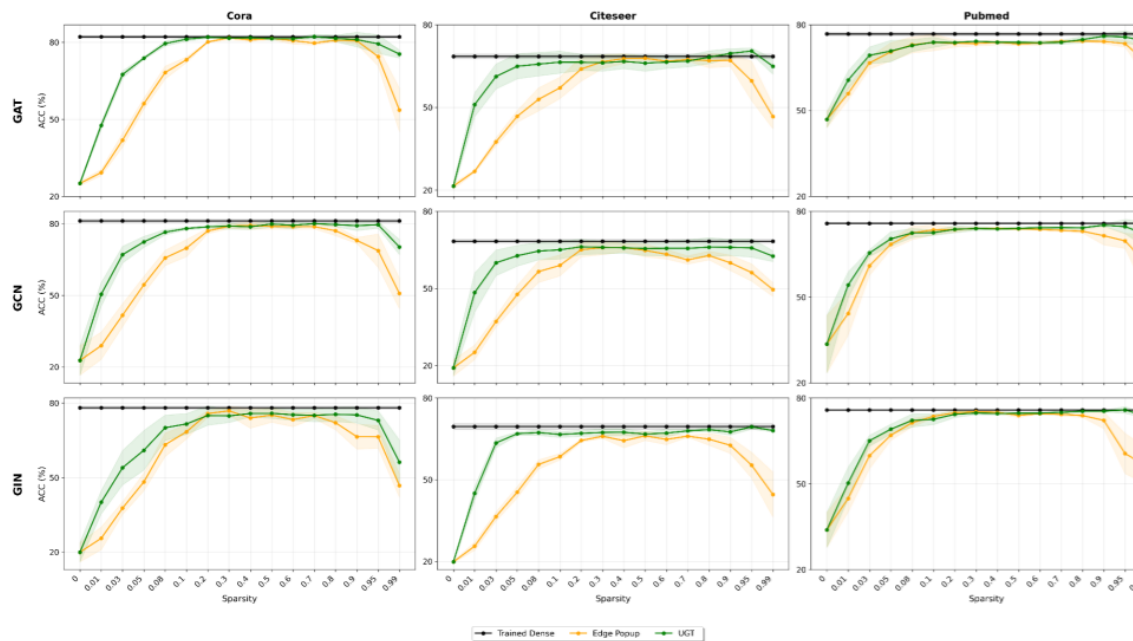


Figure 3: Test accuracy as a function of sparsity across architectures and datasets. UGT maintains competitive performance across a wide sparsity range, while Edge-Popout degrades earlier. Performance declines gradually only at extreme sparsity levels (0.95–0.99). Reproduction of Figure 5 from (1).

167 and optimizes:

$$168 \quad \min_{m_W, m_A} L(f(X, A \odot m_A; W \odot m_W), y) \quad \text{s.t.} \quad \|m_W\|_0 \leq s_W, \|m_A\|_0 \leq s_A. \quad (13)$$

169 Here, $\|\cdot\|_0$ denotes the ℓ_0 pseudo-norm, which counts the number of nonzero entries in the corresponding binary mask.

171 After training, hard pruning is applied.

172 4.1.3 Experimental Setup

173 We compare: trained dense GNN, trained sparse GNN via UGS, untrained sparse GNN (UGTs).

174 Depth is varied while holding dataset, backbone architecture, and training protocol fixed.

175 4.1.4 Findings

176 We observe that UGS collapses rapidly as depth increases, similar to trained dense GNNs. Although sparsity is enforced, representation collapse still occurs.

178 In contrast, UGTs remains stable across increasing depth.

179 This demonstrates that sparsity alone does not prevent oversmoothing under training. Instead, weight optimization appears to amplify feature homogenization, even when the model is sparse.

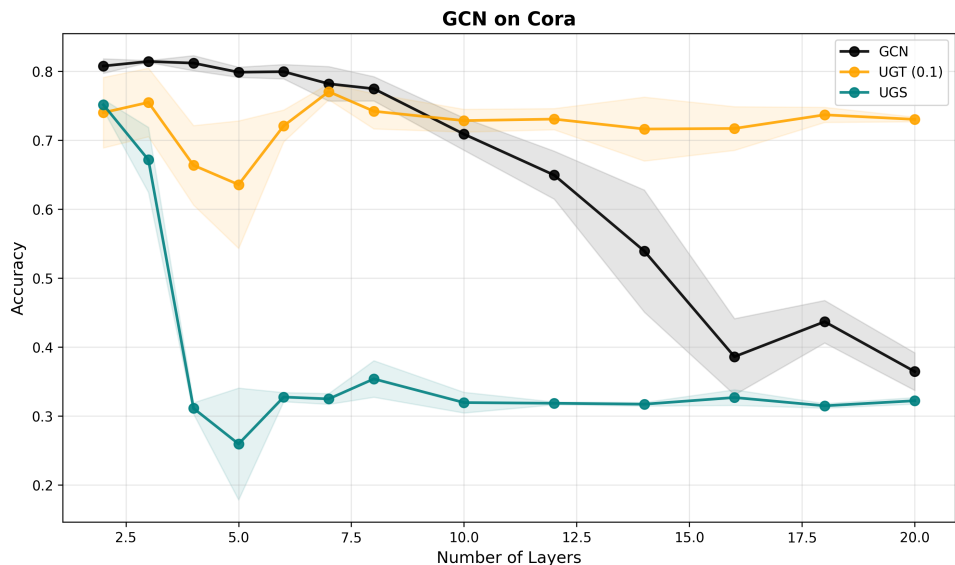


Figure 4: Accuracy versus depth comparing trained dense GNNs, UGS, and UGTs.

181 4.2 Weight Reparameterization and Initialization

182 4.2.1 Motivation

183 Recent work (5) suggests that oversmoothing is not inevitable if weights are constrained to pre-
 184 serve representation diversity. WeightRep dynamically reparameterizes weights to maintain feature
 185 correlation structure during training.

186 We compare: standard dense training, WeightInit (structured initialization only), WeightRep (dy-
 187 namic reparameterization).

188 4.2.2 WeightRep Formulation

189 Let the pre-activation features at layer ℓ be:

$$190 \hat{H}^{(\ell)} = PH^{(\ell-1)}W^{(\ell)}, \quad (14)$$

191 where P is the normalized propagation matrix.

192 Define the covariance (Gram) matrix:

$$193 \Sigma^{(\ell)} = \hat{H}^{(\ell)\top} \hat{H}^{(\ell)} + \varepsilon I, \quad (15)$$

194 where $\varepsilon > 0$ ensures numerical stability.

195 Let $f(\cdot)$ be a spectral function applied to eigenvalues. A common choice is whitening with $f(\lambda) =$
 196 $\lambda^{-1/2}$.

197 The constructed transformation matrix is:

$$\hat{W}^{(\ell)} = \left(\Sigma^{(\ell)} \right)^f. \quad (16)$$

The updated representation is then:

$$H^{(\ell)} = \hat{H}^{(\ell)} \hat{W}^{(\ell)}. \quad (17)$$

Equivalently, this induces an effective weight matrix:

$$\tilde{W}^{(\ell)} = W^{(\ell)} \hat{W}^{(\ell)}, \quad (18)$$

so that:

$$H^{(\ell)} = PH^{(\ell-1)} \tilde{W}^{(\ell)}. \quad (19)$$

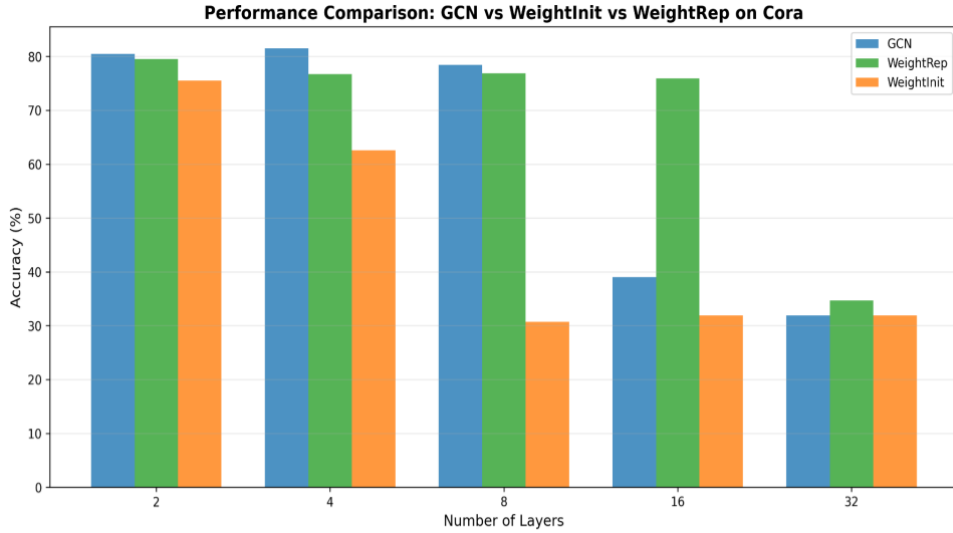


Figure 5: Depth performance comparison between dense training, WeightInit, and WeightRep.

4.3 WeightInit Formulation

Instead of modifying weights during training, we apply a one-time representation adjustment at initialization.

Let $W^{(\ell)}$ be initialized using a standard random initialization (e.g., Glorot initialization).

After computing the initial propagated features $\hat{H}^{(\ell)} = PH^{(\ell-1)}W^{(\ell)}$, we construct:

$$\Sigma^{(\ell)} = \hat{H}^{(\ell)\top} \hat{H}^{(\ell)} + \varepsilon I, \quad (20)$$

and define:

$$\hat{W}^{(\ell)} = \left(\Sigma^{(\ell)} \right)^f. \quad (21)$$

213 The initialized representation becomes:

$$214 \quad H^{(\ell)} = \hat{H}^{(\ell)} \hat{W}^{(\ell)}. \quad (22)$$

215 **4.3.1 Findings**

216 Standard dense training degrades rapidly with depth. WeightInit provides limited improvement,
 217 indicating that careful initialization alone cannot prevent collapse.

218 WeightRep, however, maintains strong performance at larger depths.

219 This suggests that oversmoothing arises from training dynamics rather than purely architectural
 220 limitations. Proper weight evolution during training is critical.

221 **4.4 UGT Initialization Study**

222 Finally, we investigate whether UGT robustness depends on random initialization.

223 We evaluate UGT under standard Kaiming initialization and WeightInit-style structured initial-
 224 ization.

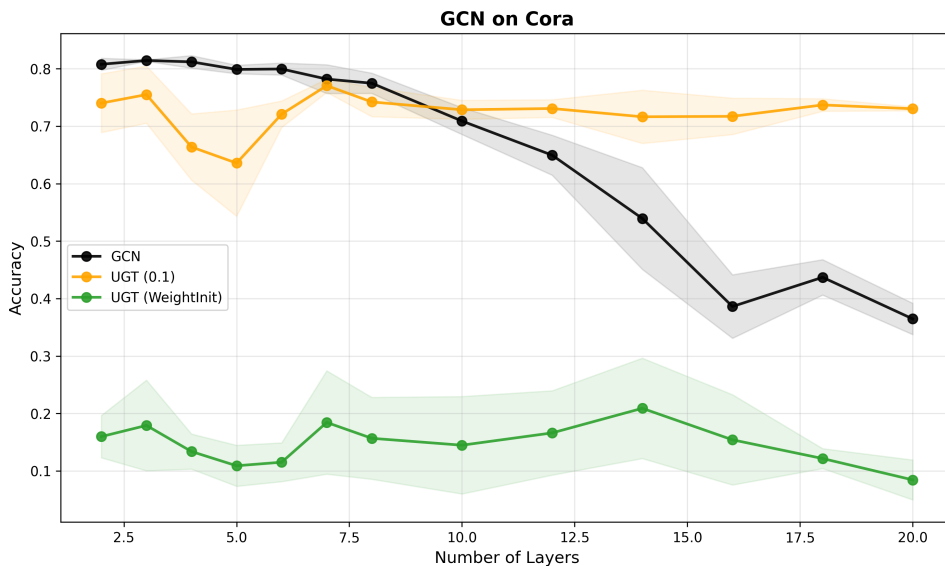


Figure 6: UGT performance under different initialization strategies.

225 **4.4.1 Findings**

226 Standard UGT remains robust across depth. Surprisingly, applying WeightInit severely degrades
 227 UGT performance.

228 This indicates that UGT robustness relies on high-entropy, unstructured random weights. Struc-
 229 tured initialization constrains the representational space in a way that conflicts with mask-based
 230 selection.

231 Thus, UGT and WeightRep prevent oversmoothing through fundamentally different mechanisms:
 232 UGT avoids destructive training dynamics entirely, while WeightRep modifies them.

233 4.5 Synthesis of Extension Results

234 Across all experiments, a consistent pattern emerges:

- 235 • Dense training amplifies representation collapse.
- 236 • Sparsity without training (UGTs) mitigates collapse.
- 237 • Sparsity with training (UGS) does not prevent collapse.
- 238 • Structured dynamic weights (WeightRep) can preserve diversity under training.

239 These results suggest that oversmoothing is primarily a learning dynamics failure rather than a
 240 structural inevitability of message passing.

241 Controlling representation evolution — either by eliminating weight training (UGTs) or constrain-
 242 ing it dynamically (WeightRep) — is key to enabling deep GNNs.

243 5 Discussion

244 Our results collectively suggest that oversmoothing is not an unavoidable consequence of depth
 245 in Graph Neural Networks, but rather a byproduct of how representations evolve under training
 246 dynamics.

247 5.1 Is Oversmoothing Structural or Learned?

248 Traditional explanations attribute oversmoothing to repeated multiplication by graph propagation
 249 operators, which act as low-pass filters. While this spectral perspective explains why representa-
 250 tions contract with depth, our experiments reveal that structural smoothing alone does not fully
 251 account for collapse.

252 Specifically:

- 253 • UGTs — which use the same propagation operators as dense GNNs — do not collapse at
 254 comparable depths.
- 255 • UGS — which enforces sparsity but retains training — still collapses.

256 Thus, sparsity of connectivity is insufficient on its own. The critical factor appears to be *how*
 257 *weights are optimized*.

258 5.2 The Role of Training Dynamics

259 Training dense weights introduces strong coupling between layers. Gradients encourage homoge-
 260 nization of features across neighborhoods, accelerating the contraction of representation space.

261 UGTs avoid this phenomenon entirely by freezing weights. Instead of adapting weights, they adapt
 262 connectivity via mask learning. This decouples representation propagation from destructive weight
 263 updates.

264 WeightRep provides further evidence that training can be compatible with depth — provided
 265 weight evolution is constrained to preserve feature diversity. Unlike UGTs, which eliminate weight
 266 training, WeightRep modifies the geometry of weight updates.

267 **5.3 Entropy and Initialization**

268 Our UGT initialization study reveals a subtle but important insight: robustness depends on high-
269 entropy, unstructured weights. When structured initialization (WeightInit) is imposed, UGT per-
270 formance degrades.

271 This suggests that mask learning in UGT relies on a rich random feature basis from which useful
272 subnetworks can be selected. Constraining this basis reduces expressivity.

273 **6 Conclusion**

274 In this work, we reproduced the principal findings of Untrained Graph Neural Networks Tickets
275 and systematically investigated the mechanisms underlying oversmoothing in deep GNNs.

276 Our reproduction confirms three central observations:

- 277 1. Standard dense GNNs suffer severe performance degradation beyond moderate depth.
- 278 2. Untrained sparse subnetworks (UGTs) maintain stable accuracy even at 20+ layers.
- 279 3. UGTs preserves representation diversity, as measured by Mean Average Distance (MAD),
280 whereas dense models rapidly collapse.

281 Building on this foundation, our extensions clarify the role of sparsity, initialization, and train-
282 ing:

- 283 • Sparsity alone does not prevent oversmoothing, as demonstrated by the rapid collapse of
284 UGS under training.
- 285 • Proper weight evolution (WeightRep) enables deep training without collapse, indicating that
286 oversmoothing is not structurally inevitable.
- 287 • UGT robustness depends on high-entropy random initialization, highlighting the importance
288 of expressive random feature bases.

289 Collectively, these findings support a central thesis:

290 Oversmoothing is primarily a learning dynamics failure rather than a fundamental
291 architectural limitation of message passing.

292 Deep GNN failure arises when training drives representations toward a low-rank, homogeneous
293 subspace. Preventing collapse requires either eliminating destructive weight updates (UGTs) or
294 actively constraining weight evolution (WeightRep).

295 Future work may explore hybrid approaches that combine sparse connectivity, dynamic weight
296 control, and theoretical guarantees on representation rank preservation. Understanding how rep-
297 resentation geometry evolves during training remains a promising direction for building stable,
298 deep graph models.

299 **References**

- 300 [1] T. Huang et al. You Can Have Better Graph Neural Networks by Not Training Weights at All:
301 Finding Untrained GNNs Tickets. In Learning on Graphs Conference (LoG), 2022.
- 302 [2] V. Ramanujan et al. What’s Hidden in a Randomly Weighted Neural Network? In CVPR,
303 2020.

- 304 [3] D. Chen et al. Measuring and Relieving the Over-smoothing Problem for Graph Neural Net-
305 works from the Topological View. In AAAI, 2020.
- 306 [4] T. Chen et al. A Unified Lottery Ticket Hypothesis for Graph Neural Networks. In ICML,
307 2021.
- 308 [5] Z. Zhuo et al. Graph Neural Networks (with Proper Weights) Can Escape Oversmoothing. In
309 ACML, 2024.