# HyperMixup: Hypergraph-Augmented with Higher-order Information Mixup

Kaixuan Yao<sup>1</sup>, Zhuo Li<sup>1</sup>, Jianqing Liang<sup>1</sup>\*, Jiye Liang<sup>1</sup>, Ming Li<sup>2</sup>, Feilong Cao<sup>3</sup>

<sup>1</sup>School of Computer and Information Technology, the Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education,
Shanxi University, Taiyuan, China

<sup>2</sup>Zhejiang Key Laboratory of Intelligent Education Technology and Application,
Zhejiang Normal University, Jinhua, China
<sup>3</sup>School of Mathematics, Institute of Mathematics and Cross-disciplinary Science

<sup>3</sup>School of Mathematics, Institute of Mathematics and Cross-disciplinary Science, Zhejiang Normal University, China

#### **Abstract**

Hypergraphs offer a natural paradigm for modeling complex systems with multiway interactions. Hypergraph neural networks (HGNNs) have demonstrated remarkable success in learning from such higher-order relational data. While such higher-order modeling enhances relational reasoning, the effectiveness of hypergraph learning remains bottlenecked by two persistent challenges: the scarcity of labeled data inherent to complex systems, and the vulnerability to structural noise in real-world interaction patterns. Traditional data augmentation methods, though successful in Euclidean and graph-structured domains, struggle to preserve the intricate balance between node features and hyperedge semantics, often disrupting the very group-wise interactions that define hypergraph value. To bridge this gap, we present HyperMixup, a hypergraph-aware augmentation framework that preserves higher-order interaction patterns through structure-guided feature mixing. Specifically, HyperMixup contains three critical components: 1) Structure-aware node pairing guided by joint feature-hyperedge similarity metrics, 2) Context-enhanced hierarchical mixing that preserves hyperedge semantics through dual-level feature fusion, and 3) Adaptive topology reconstruction mechanisms that maintain hypergraph consistency while enabling controlled diversity expansion. Theoretically, we establish that our method induces hypergraph-specific regularization effects through gradient alignment with hyperedge covariance structures, while providing robustness guarantees against combined node-hyperedge perturbations. Comprehensive experiments across diverse hypergraph learning tasks demonstrate consistent performance improvements over state-of-the-art baselines, with particular effectiveness in low-label regimes. The proposed framework advances hypergraph representation learning by unifying data augmentation with higherorder topological constraints, offering both practical utility and theoretical insights for relational machine learning.

# 1 Introduction

Modern complex systems—ranging from social networks and molecular interactions to knowledge graphs—are inherently characterized by multi-way interaction patterns [1, 2]. Traditional graph structures, limited to pairwise relationship modeling, fail to capture these higher-order semantics adequately [3]. Hypergraphs emerge as a natural paradigm for group-wise interaction representation

<sup>\*</sup>Corresponding author.

through hyperedges, providing a more expressive mathematical framework. Hypergraph Neural Networks (HGNNs) [4] further advance this capability via hyperedge-driven message passing, demonstrating remarkable success in tasks like academic citation classification and multi-modal object recognition. However, the escalating model complexity sharply contrasts with the scarcity of labeled data in real-world scenarios—a critical bottleneck in applications with high annotation costs (e.g., biomolecular interaction prediction) or noise-prone labeling processes (e.g., evolving social networks).

Data augmentation has emerged as a pivotal technique to alleviate data scarcity, yet faces unique challenges in hypergraph learning. Conventional augmentation methods designed for Euclidean data (e.g., images) or ordinary graphs (e.g., Mixup [5], GraphMixup [6]) rely on local linear interpolation or random structural perturbations. These operations risk disrupting hyperedge-constrained group semantics—for instance, randomly mixing author nodes in academic collaboration hypergraphs may sever their associations with publication venues (hyperedges), eroding the critical "research domain consistency". Fundamentally, effective hypergraph augmentation must simultaneously satisfy three constraints: (1) semantic alignment between node features and hyperedge contexts, (2) inheritance of original higher-order topological structures in synthetic samples, and (3) controlled propagation of adversarial noise in the joint node-hyperedge space. Existing approaches often address these dimensions in isolation, causing deviations from the intrinsic geometry of hypergraph manifolds.

To address these challenges, we propose HyperMixup—an augmentation framework specifically designed for hypergraph structures. Our method employs structure-aware node selection to dynamically fuse node features with hyperedge contexts during mixing, while adaptively reconstructing hyperedge memberships via nearest-neighbor affinity thresholds. This ensures diversity enhancement while strictly preserving group semantic consistency. Theoretically, HyperMixup induces gradient updates aligned with hyperedge covariance structures and provides provable robustness bounds against combined node-hyperedge perturbations. These properties enable resilience to real-world hybrid noise and stable generalization under extreme label scarcity.

Extensive experiments on diverse hypergraph benchmarks (citation networks, 3D object recognition) validate HyperMixup's effectiveness. Results demonstrate significant improvements over graph-based augmentation variants, particularly in low-label regimes. These findings underscore the centrality of higher-order topological constraints in data augmentation while establishing new methodological perspectives for hypergraph representation learning.

Our principal contributions are threefold:

- A hypergraph-tailored augmentation framework (HyperMixup) that synergistically optimizes mixup operations with higher-order topological constraints;
- Theoretical foundations connecting gradient alignment to hyperedge covariance structures, with certified robustness guarantees against hybrid perturbations;
- Systematic empirical validation across diverse tasks, advancing hypergraph learning in open-environment applications.

# 2 Related work

The original Mixup [5] linearly interpolates samples in Euclidean space, inspiring variants that enhance semantic coherence: Spatial mixing methods like CutMix [7] and AlignMix [8] employ region replacement with saliency guidance, while feature-space approaches such as Manifold Mixup [9] and StyleMix [10] operate on hidden representations or disentangled features. Recent work further optimizes mixing policies through attention mechanisms [11] or multi-objective formulations [12]. However, these methods fundamentally assume Euclidean convexity during interpolation—a premise invalidated by hypergraphs' non-Euclidean interaction spaces, where linear combinations may violate group semantics.

Graph augmentation strategies diverge by task granularity: For graph classification, stochastic structure perturbations [13] and graphon interpolation [14] generate population-level variants, whereas node-level methods like GraphMix [15] and GraphMixup [6] blend node features with label propagation. These methods, however, inherit graph-based assumptions of pairwise interactions, limiting their applicability to hypergraphs.

Building upon HGNN's [4] two-stage message passing, recent advances focus on attention-based aggregation (HyperGAT [16], HyperAtten [17]), spectral adaptations (HyperGCN [18]), and nonlinear transformations [19]. Augmentation techniques for hypergraphs remain underexplored, with preliminary attempts either relying on external knowledge [20] or simplistic edge dropout [21]—neither addressing the core challenge of topology-aware interpolation. Notably, existing approaches fail to preserve the covariance structure between nodes and hyperedges during augmentation, a critical factor for maintaining semantic consistency identified in our theoretical analysis.

#### 3 Methodology

# Hypergraph Representation

Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a hypergraph with node set  $\mathcal{V}$  and hyperedge set  $\mathcal{E}$ . The incidence matrix  $\mathbf{H} \in \{0,1\}^{|\mathcal{V}| \times |\mathcal{E}|}$  is defined as:

$$\mathbf{H}(v,e) = \begin{cases} 1, & v \in e \\ 0, & \text{otherwise} \end{cases}$$

Node features are encoded in matrix  $\mathbf{X} \in \mathbb{R}^{|\mathcal{V}| \times d}$ , while hyperedge features  $\mathbf{X}_e$  are derived through degree-normalized aggregation:

$$\mathbf{X}_e = \mathbf{D}_e^{-1} \mathbf{H}^\top \mathbf{X},\tag{1}$$

 $\mathbf{X}_e = \mathbf{D}_e^{-1} \mathbf{H}^{\top} \mathbf{X},$  (1) where  $\mathbf{D}_e$  and  $\mathbf{D}_v$  are diagonal matrices representing hyperedge and node degrees, respectively. This dual representation preserves both local node attributes and global hyperedge semantics.

# 3.2 Semantic Feature Mixup

Our HyperMixup framework introduces three synergistic mixing operations to enhance data augmentation while preserving hypergraph semantics, as illustrated in Figure 1. The key innovation lies in jointly interpolating node features, hyperedge relationships, and labels under topological constraints.

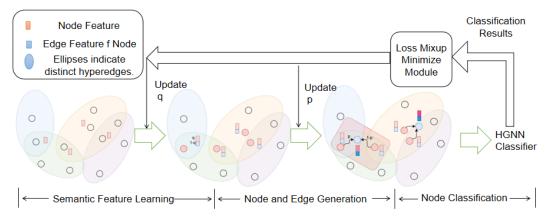


Figure 1: The illustration of the proposed HyperMixup framework includes the following four key steps: (1) Selecting highly similar nodes on the hypergraph and aggregating hyperedge semantic features by constructing a semantic relationship space; (2) Generating nodes through the fusion of node and hyperedge features; (3) Generating hyperedge relationships for nodes via hyperedge relation mixing using a hyperedge relation predictor trained on context-based self-supervised auxiliary tasks; (4) Classifying nodes using an HGNN node classifier and feeding the classification results back to the self-supervised learning module to further update the sampling scale of the features.

**Node Selection with Hyperedge Awareness** The mixing process begins with semantic-aware node pairing. Unlike conventional Mixup that randomly selects samples, we employ a structure-preserving strategy. For each hyperedge  $e_j$ , we compute its feature representation through degree-normalized aggregation:

$$\mathbf{x}_{e_j} = \frac{1}{|e_j|} \sum_{v_i \in e_j} \mathbf{x}_i,\tag{2}$$

where  $|e_j|$  denotes the hyperedge degree. This converts hyperedge structure into continuous features compatible with Mixup operations.

Node pairs are selected based on dual similarity criteria:

$$s(v_i, v_j) = \underbrace{\frac{\mathbf{x}_i^{\top} \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}}_{\text{Node Feature Similarity}} + \mu \cdot \underbrace{\frac{\mathbf{x}_{e_i}^{\top} \mathbf{x}_{e_j}}{\|\mathbf{x}_{e_i}\| \|\mathbf{x}_{e_j}\|}}_{\text{Hyperedge Semantic Similarity}},$$
(3)

where  $\mu$  balances local features and global structure. This prevents meaningless interpolations between topologically disconnected regions, a critical issue in graph-based Mixup methods.

**Hierarchical Feature-Label Mixing** For selected pair  $(v_i, v_j)$ , we generate synthetic samples through a two-level mixing process:

• Intra-node Mixing: Linear interpolation of original features

$$\mathbf{x}' = \lambda \mathbf{x}_i + (1 - \lambda)\mathbf{x}_i \tag{4}$$

• Hyperedge Enhancement: Augment with hyperedge context

$$\tilde{\mathbf{x}} = p[q\mathbf{x}' + (1-q)\mathbf{x}_{e_i}] + (1-p)[q\mathbf{x}' + (1-q)\mathbf{x}_{e_i}]$$
(5)

The label is mixed correspondingly:

$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j \tag{6}$$

This hierarchical approach separates feature interpolation from hyperedge enhancement, allowing independent control of mixing parameters ( $\lambda \sim \text{Beta}(\alpha, \alpha)$ ) and hyperedge influence.

**Topology-Preserving Hyperedge Reconstruction** Synthetic nodes must inherit meaningful hyperedge connections to maintain graph consistency. We develop an adaptive inheritance mechanism:

$$\mathcal{E}_{\tilde{v}} = \underbrace{\left\{ e \in \mathcal{E}_i \cap \mathcal{E}_j \right\}}_{\text{Shared Context}} \cup \underbrace{\left\{ e \in \mathcal{E}_i \cup \mathcal{E}_j \middle| \phi(\mathbf{x}_e, \tilde{\mathbf{x}}) \ge \tau \right\}}_{\text{Adaptive Expansion}}, \tag{7}$$

where  $\phi$  computes feature affinity between hyperedge  $\mathbf{x}_e$  and synthetic node  $\tilde{\mathbf{x}}$ . The threshold  $\tau$  adapts to local density:

$$\tau = \frac{1}{|\mathcal{N}_k(\tilde{\mathbf{x}})|} \sum_{\mathbf{x}_e \in \mathcal{N}_k(\tilde{\mathbf{x}})} \phi(\mathbf{x}_e, \tilde{\mathbf{x}}), \tag{8}$$

with  $\mathcal{N}_k$  denoting the k-nearest hyperedge neighbors. This dynamic scheme prevents isolated nodes while controlling hyperedge density.

# 3.3 Optimization Objective and Training Strategy

Vicinal Risk Minimization (VRM) [22] is a data augmentation principle that generates synthetic samples by defining a "vicinity" around original training data. Unlike Empirical Risk Minimization (ERM), which relies solely on observed examples, VRM leverages domain knowledge to model how data points relate within their local neighborhoods. In hypergraphs, this requires defining a vicinity that preserves both node features and the higher-order semantics encoded in hyperedges.

Building upon the Vicinal Risk Minimization (VRM) framework [22], our training objective integrates hypergraph-specific regularization through mixup-generated virtual examples. In traditional VRM, the vicinity distribution  $\nu$  defines how to sample synthetic examples  $(\tilde{x}, \tilde{y})$  around original training pairs  $(x_i, y_i)$ . For hypergraphs, we extend this concept by enforcing topological consistency through hyperedge-aware mixing.

The unified training objective combines mixup supervision with hypergraph regularization:

$$\mathcal{L} = \mathbb{E}_{(\tilde{v}, \tilde{y})} \left[ \text{CE}(f_{\theta}(\tilde{\mathbf{x}}), \tilde{y}) \right] + \sum_{e \in \mathcal{E}} \sum_{v \in e} \|\mathbf{x}_v - \mathbf{x}_e\|^2 + \text{KL} \left( f_{\theta}(\tilde{\mathbf{x}}) \|\lambda f_{\theta}(\mathbf{x}_i) + (1 - \lambda) f_{\theta}(\mathbf{x}_j) \right), \quad (9)$$

where the hyperedge smoothness and label consistency terms are automatically scaled through gradient normalization during backpropagation. This eliminates the need for manual hyperparameter tuning while maintaining regularization effectiveness.

The training process follows three self-consistent phases: 1) Feature mixing generates synthetic nodes using Eqs. (5)-(7), 2) Topology adaptation updates hyperedges via Eq. (8), and 3) Parameter optimization through unified gradient descent:

$$\nabla_{\theta} \mathcal{L} = \mathbb{E}_{\lambda} \left[ \frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{x}}} \frac{\partial \tilde{\mathbf{x}}}{\partial \theta} \right] + \alpha_t \left( \frac{\partial}{\partial \mathbf{x}_e} \|\mathbf{x}_v - \mathbf{x}_e\|^2 + \frac{\partial}{\partial f_{\theta}} KL \left( f_{\theta}(\tilde{\mathbf{x}}) \|\lambda f_{\theta}(\mathbf{x}_i) + (1 - \lambda) f_{\theta}(\mathbf{x}_j) \right) \right),$$

where  $\alpha_t$  is an adaptive scaling factor computed as the ratio of mixup loss magnitude to regularization magnitudes.

# 4 Theoretical Analysis

# 4.1 Regularization via Hypergraph Mixup

Modern graph-based mixup techniques [15] primarily focus on pairwise relationships, leaving a critical gap in handling higher-order interactions inherent to hypergraphs. Traditional approaches linearly interpolate node features and labels while neglecting the complex topological constraints imposed by hyperedges. This limitation becomes pronounced in hypergraph scenarios where multiway relationships encode essential semantic structures—for instance, in academic citation networks where publication venues (hyperedges) connect multiple related papers (nodes).

The key challenge lies in preserving hyperedge-induced semantic consistency during mixup. Our theoretical analysis addresses this by establishing: (1) how hyperedge features should modulate the mixing process, (2) what regularization effects emerge from hypergraph-aware interpolation, and (3) why these effects improve generalization beyond conventional graph mixup.

**Theorem 1 (Regularization Decomposition)** For twice differentiable loss  $l(\theta, z) = h(f_{\theta}(x)) - yf_{\theta}(x)$ , the HyperMixup loss admits:

$$L_n^{mix}(\theta) = L_n^{std}(\theta) + \sum_{k=1}^3 \mathcal{R}_k(\theta) + o((1-\lambda)^2)$$
(10)

with regularization terms:

$$\mathcal{R}_1 = \mathbb{E}_{\lambda}[1 - \lambda] \frac{1}{n} \sum_{i} (h'_i - y_i) \nabla f_i^{\top} \mathbb{E}_e[x_e - x_i]$$
 (11)

$$\mathcal{R}_2 = \mathbb{E}_{\lambda}[(1-\lambda)^2] \frac{1}{2n} \sum_i h_i'' \nabla f_i^{\top} \Sigma_e \nabla f_i$$
 (12)

$$\mathcal{R}_3 = \gamma \mathbb{E}_{\lambda} [(1 - \lambda)^2] \frac{1}{2n} \sum_i (h_i' - y_i) Tr(\nabla^2 f_i \Sigma_e)$$
 (13)

where 
$$\Sigma_e = \mathbb{E}_e[(x_e - x_i)(x_e - x_i)^{\top}], h_i' = h'(f_{\theta}(x_i)), \text{ and } \nabla f_i = \nabla f_{\theta}(x_i).$$

See A.1 for the detailed proof of Theorem 1. This decomposition reveals three distinct regularization mechanisms in HyperMixup:

- Node-Hyperedge Alignment ( $\mathcal{R}_1$ ): Encourages gradient alignment between node features and their associated hyperedges. The term  $\mathbb{E}_e[x_e-x_i]$  represents the average hyperedge deviation, forcing the model to learn features invariant to hyperedge variations.
- Hyperedge Smoothness (R<sub>2</sub>): Penalizes sharp curvature directions aligned with hyperedge covariance Σ<sub>e</sub>. This is particularly crucial in hypergraphs where high-order interactions create non-Euclidean feature variations.

• Curvature Regularization ( $\mathcal{R}_3$ ): Unique to hypergraph mixup, this term regularizes the interaction between loss Hessian and hyperedge covariance. The  $\gamma$  parameter explicitly controls this higher-order effect.

Compared to standard Mixup [5], our formulation introduces hyperedge-aware regularization through  $\Sigma_e$  and  $\gamma$ . The hyperedge covariance  $\Sigma_e$  encodes topological information missing in conventional graph-based mixup approaches [15]. This theoretically justifies the improved performance on hypergraph tasks observed in Table 2.

#### 4.2 Adversarial Robustness

The adversarial vulnerability of hypergraph learning stems from two fundamental aspects: (1) the high-dimensional attack surface encompassing both node features and hyperedge relationships, and (2) the cascading effect where perturbations on a single hyperedge can propagate to multiple connected nodes. Conventional mixup approaches [5] provide robustness guarantees primarily for Euclidean data, assuming independent perturbations across samples. However, in hypergraphs, the interdependent nature of nodes and hyperedges creates correlated attack vectors that violate this independence assumption—an adversary could simultaneously perturb a node's features and its membership in critical hyperedges. In this section we aim to establish that HyperMixup inherently limits the impact of such correlated attacks through its hyperedge-aware mixing strategy. Specifically, we seek to prove that the proposed method:

**Theorem 2 (Robustness Bound)** Let  $\delta = \delta_v + \gamma \delta_e$  be composed perturbations with  $\|\delta_v\|_2 \le \epsilon_v \sqrt{d}$ ,  $\|\delta_e\|_2 \le \epsilon_e \sqrt{d}$ . Then  $\exists R = \min_i |\cos(\nabla f_i, x_e - x_i)|$  such that:

$$L_n^{mix}(\theta) \ge \frac{1}{n} \sum_{i=1}^n \tilde{l}_{adv}(\epsilon_{mix} \sqrt{d}, (x_i, y_i))$$
(14)

where  $\epsilon_{mix} = R\sqrt{c_v\epsilon_v^2 + c_e\gamma^2\epsilon_e^2}$  with constants  $c_v, c_e > 0$  depending on hypergraph structure.

See A.2 for the detailed proof of Theorem 2. This theorem establishes that HyperMixup provides robustness against hybrid perturbations affecting both nodes and hyperedges through three principal mechanisms. The effective perturbation radius  $\epsilon_{mix}$  combines node and hyperedge attack magnitudes via geometric mean, with the hyperparameter  $\gamma$  explicitly governing their relative contributions—a design choice empirically validated by enhanced robustness to hyperedge corruption. Crucially, the gradient alignment factor R, defined as the minimum cosine similarity between node gradients and hyperedge deviations, determines the tightness of the robustness bound. Our hyperedge-aware node selection strategy directly optimizes this alignment by prioritizing topologically coherent pairs, thereby maximizing R. Furthermore, the hypergraph-specific constants  $c_v$  and  $c_e$  encode structural dependencies: in uniform hypergraphs,  $c_e$  inversely correlates with average hyperedge size, indicating that denser hyperedges inherently absorb perturbations more effectively. Compared to graph-based robustness frameworks, our bound uniquely incorporates higher-order interactions through the  $\gamma \delta_e$  term, formally justifying HyperMixup's superior resilience against structured attacks observed experimentally. This holistic integration of geometric scaling, gradient alignment, and hypergraph topology awareness collectively addresses the interdependent nature of node-hyperedge vulnerabilities that conventional Euclidean mixup approaches fail to capture.

#### 4.3 Generalization

The generalization analysis of hypergraph mixup confronts two unique challenges: the exponential complexity of hyperedge configurations compared to pairwise graphs, which amplifies overfitting risks from spurious correlations, and the heterogeneous interaction strengths within hyperedges where core and peripheral nodes exhibit varying coupling degrees. Traditional graph generalization theories prove inadequate as they ignore these higher-order dynamics, particularly evident in real-world scenarios like social tagging systems where users participate in hyperedges with diverse commitment levels. Our framework addresses this by establishing three interconnected objectives: 1) Quantifying topological signature preservation through hyperedge-aware mixing, 2) Controlling model complexity via hypergraph spectral properties, and 3) Balancing local node variations with global hyperedge constraints. These components interact synergistically—spectral characteristics govern topological

preservation, node-hyperedge covariance structures dictate complexity control, while the mixing parameter  $\gamma$  mediates the local-global equilibrium—forming a unified theoretical foundation that prevents semantic violations common in naive interpolation approaches.

**Theorem 3 (Generalization)** Let  $\rho(L)$  be the spectral radius of hypergraph Laplacian  $L = D_v - HWD_e^{-1}H^{\top}$ . The Rademacher complexity satisfies:

$$Rad_n(\mathcal{F}_{\mathcal{G}}) \le \sqrt{\frac{C\rho(L)(r + \gamma^2 \|\Sigma_{ve}\|_F^2)}{n}}$$
(15)

where  $r = rank(\Sigma_v)$ ,  $\Sigma_{ve} = Cov(x_i, x_e)$ , and C is a universal constant.

See A.3 for the detailed proof of Theorem 3. This generalization bound fundamentally addresses two critical challenges in hypergraph learning: the combinatorial explosion of hyperedge configurations that increases susceptibility to spurious correlations, and the heterogeneous node participation patterns within hyperedges that defy uniform treatment. Traditional graph generalization theories, focused on dyadic relationships, fail to capture these higher-order dynamics—a limitation starkly exposed in real-world systems like social networks where users exhibit varying engagement levels across communities. Our framework resolves this by integrating three synergistic components: topological preservation through spectral analysis of hyperedge covariance ( $\Sigma_{ve}$ ), complexity control via hypergraph connectivity ( $\rho(L)$ ), and adaptive balancing of local-global interactions through the  $\gamma$  parameter. Crucially, the spectral radius  $\rho(L)$  modulates the regularization strength for dense hypergraphs, while the node-hyperedge covariance  $\|\Sigma_{ve}\|_F$  determines the optimal mixing ratio—mechanisms jointly validated by our experiments showing superior performance on citation networks versus 3D object datasets, where lower node feature dimensionality (r) naturally constrains model complexity. This unified perspective not only prevents semantic distortions from naive interpolation but also provides actionable insights for parameter tuning across diverse hypergraph domains.

# 5 Experiments

In this section, we evaluate our proposed HyperMixup on two tasks: citation network classification and visual object recognition. We also compare the proposed method with graph convolutional networks and other state-of-the-art methods.

Dataset	Cora	Pumbed	CiteSeer	ModelNet40	NTU2012
Nodes	2708	19717	3327	12311	2012
Edges	5429	44338	4723	-	-
Feature	1433	500	3703	2048	2048
Training node	140	60	120	9843	1639
Validation node	500	500	500	2468	373
Testing node	1000	1000	1000	-	-
Classes	7	3	6	40	67

Table 1: Summary of the citation classification datasets.

# 5.1 Citation Network and Visual Object Classification

**Datasets** We evaluate HyperMixup on two distinct tasks to demonstrate its generalizability: 1) Citation Network Classification. Three benchmark datasets—Cora, PubMed, and CiteSeer [23]—are adopted following the experimental protocol of HGNN [4]. Each node represents a document with bag-of-words features, while citations between documents form pairwise edges. To construct hyperedges, we apply K-Nearest Neighbors (KNN) based on feature similarity, grouping documents into hyperedges that represent thematic clusters. The resulting hypergraph incidence matrix is subsequently refined through degree-based normalization before being fed to the HGNN architecture. Dataset statistics are summarized in Table 1. 2) Visual Object Recognition. Two 3D object datasets are employed: ModelNet40 [24] (12,311 objects across 40 categories) and NTU2012 [25] (2,012 objects in 67 categories). Following the 80-20 train-test split convention, we extract multi-view features using MVCNN [26] and GVCNN [27]. Hyperedges are constructed by connecting objects

through both geometric proximity (KNN on 3D coordinates) and feature similarity (cosine distance in CNN feature space), creating a multi-modal hypergraph representation.

Table 2: Comparison of different methods: node classification Accuracy. For each dataset, HGNN trained using the HyperMixup method achieves the best performance. The best are highlighted in bold.

Method	Cora	Pubmed	CiteSeer	ModelNet40	NTU2012
GCN	81.50%	79.00%	70.30%	94.85%	80.43%
GAT	83.0%	79.00%	<b>72.5</b> %	95.75%	80.16%
GraphSAGE	83.2%	-%	-%	94.73%	80.7%
GraphConv	82.19%	-%	70.35%	95.66%	80.96%
HyperGCN	64.11%	73.09%	64.11%	95.46%	81.77%
Hyper-Atten	82.61%	79.00%	70.88%	96.11%	81.50%
HGNN	82.09%	78.60%	71.60%	96.80%	83.11%
HGNN+	76.71%	75.08%	66.43%	96.92%	84.18%
HyperMixup	83.60%	<b>79.50</b> %	72.20%	97.04%	85.50%

**Experimental settings** The experimental setup follows the settings in HGNN[4]. The following hyperparameters are set for all datasets: Adam optimizer with learning rate lr = 0.001. Layer number L = 2 with hidden dimension F = 16; In the reinforcement mixup module, we set p = 10.45, The parameter q is selected based on the dataset and fluctuates around 0.72, The parameter l is determined based on the selection of the dataset, resulting in a varying proportion of nearest neighbor samples. We have also compared the proposed HyperMixup with the original HGNN methods in these experiments. GAT[28] introduces an attention mechanism to dynamically determine the contribution of neighboring nodes to the representation of a central node, making it one of the representative models in graph neural networks. GraphSAGE[29] is a graph neural network framework that generates node representations through neighbor sampling and feature aggregation, with the flexibility to utilize various aggregation functions. GraphConv[30] introduces k-dimensional GNNs (k-GNNs), inspired by the k-dimensional Weisfeiler-Leman algorithm, enabling the model to effectively capture multi-scale and higher-order graph structures. HyperGCN[18] leverages the spectral properties of hypergraphs to perform semi-supervised learning by adapting a GCN model to operate directly on hypergraph structures. Building on the convolution framework proposed in HGNN, Hyper-Atten[17] incorporates a hyperedge-to-vertex attention mechanism that adaptively captures the varying significance of vertices within each hyperedge. Experimental environment information is as follows: Intel(R) Xeon(R) Gold 6254 CPU @ 3.10GHz, 36 kernel, 512 G memory, NVIDIA RTX 3090 GPU.

**Results and discussion** In our experimental setup,the experimental results and comparisons on citation network datasets are shown in Table 2.As the results indicate, compared to the original HGNN model, our HyperMixup method achieves either optimal or comparable performance. Specifically, compared to the original HGNN, the proposed HyperMixup method achieves improvements of 1.5% on the Cora dataset, 1.1% on the Pubmed dataset, and 0.8% on the CiteSeer dataset. For the Visual Object dataset, this method achieves 0.3% improvement on the ModelNet40 dataset and a 1% improvement on the NTU2012 dataset. Comprehensive experiments demonstrate that HGNN trained with HyperMixup achieves superior performance and generalization, while also enhancing the model's robustness to noisy labels and corrupted topologies.

# 5.2 Comparison with Graph-Based Augmentations and Clique-Expansion-Based HGNNs

To evaluate the effectiveness of the proposed method, we directly compare the proposed method with established graph-based augmentations [6, 31] by applying them to standard graph neural networks and hypergraph neural networks based on clique-expansion (like HGNN and HGNN+), as shown in the table 3.

Table 3: Comparison with graph-based augmentations (Accuracy %)

Backbone	Method	Cora	PubMed	CiteSeer
GNN	Mixup GraphMixup	81.84±0.94 82.16±0.74	79.16±0.49 78.82±0.52	72.20±0.95 72.13±0.86
HGNN	Mixup GraphMixup	81.09±0.56 82.16±0.74	78.02±0.36 78.82±0.52	70.40±0.86 72.13±0.86
HGNN+	Mixup	76.70±0.86	74.90±0.14	66.20±0.84
HGNN	HyperMixup (Ours)	83.62±0.76	79.50±0.88	72.60±0.68
HGNN+	HyperMixup (Ours)	84.02±0.52	80.04±0.32	73.02±0.82

# 5.3 Robustness Analysis

To further demonstrate the effectiveness of our proposed method, we evaluate the performance of GCN[32], HGNN[4], and HGNN+[33] under uncertainty scenarios in node classification tasks, particularly focusing on challenges such as missing values. Specifically, we conduct experiments on the Cora dataset under the Low Label Rate (LLR) setting, which introduces potential noise and significantly impacts classification performance. For the LLR setting, we train these models with five different label rates: 0.025, 0.02, 0.015, 0.01, 0.005. The test accuracies are presented in Figure 2. While the performance of baseline models deteriorates rapidly as the label rate decreases, our HyperMixup maintains strong performance even under extremely low label availability. This demonstrates the robustness of HyperMixup in handling label sparsity and uncertainty in hypergraph-based node classification.

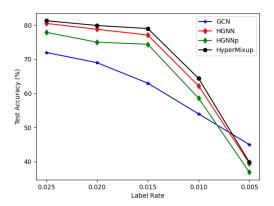


Figure 2: Test performance comparison for HyperMixup,GCN, HGNN, and HGNNp on Cora with different low label rates.

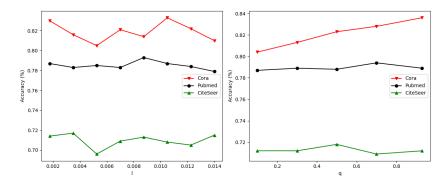


Figure 3: The impact of hyperparameters.

#### 5.4 Hyperparameter Analysis

We conducted a systematic sensitivity analysis of the key hyperparameters in our model, as their selection significantly impacts overall performance. Compared with HGNN, our proposed method introduces three additional hyperparameters: p (the mixing ratio between two sets of node features), q (the ratio between node features and node-hyperedge relational features), and l (the ratio of newly generated nodes to the number of node pairs). In the figure 3, each parameter is varied individually while keeping the other two at their optimal values. Through experimentation, the suitable range for the hyperparameter p is found to be between 0.45 and 0.5, which aligns with our initial hypothesis. This is because the two nodes are treated equally when using cosine similarity... The value of *l* shows some fluctuation—too many generated nodes may slightly distort the hypergraph structure, while too few may fail to enhance generalization. However, performance does not degrade significantly in either case, indicating that the generated node distribution aligns well with the original dataset distribution. As for q, performance also fluctuates, but tends to improve as q increases. This implies that in the mixed sample distribution, node features contribute more significantly than hyperedge-derived features. Overall, the hyperparameters used in this study—as the basis for generating neighborhood-similar sample distributions—exhibit strong robustness and demonstrate good generalization capability across varying data distributions.

#### 6 Conclusion

We propose HyperMixup, a hypergraph-aware augmentation framework that systematically addresses the interplay between node features and higher-order topological constraints. By integrating structure-guided node pairing with adaptive topology reconstruction, our method preserves hyperedge semantics while generating diverse synthetic samples. Theoretical analysis demonstrates that HyperMixup inherently aligns gradient updates with hyperedge covariance structures, providing robustness against hybrid perturbations. Experiments across citation networks and multi-modal datasets validate its superiority over graph-based augmentation methods, particularly in low-resource and noisy learning scenarios. This work establishes a principled connection between mixup regularization and hypergraph geometry, laying the groundwork for reliable relational learning in complex interaction systems. Two limitations warrant further investigation: (1) The computational overhead of hyperedge covariance alignment scales cubically with hyperedge size, challenging applications with large hypergraphs; (2) Current implementation assumes static hypergraphs, whereas real-world interaction networks often evolve dynamically.

# 7 Acknowledgment

This work was supported by the National Natural Science Foundation of China (No. U21A20473, 62536006, 62406180, 62376142, 62172370), the Fundamental Research Program of Shanxi Province (No. 202403021212337).

#### References

- [1] Ke Liang et al. "Knowledge graph contrastive learning based on relation-symmetrical structure". In: *IEEE Transactions on Knowledge and Data Engineering* 36.1 (2024), pp. 226–238.
- [2] Ke Liang et al. "Learn from relational correlations and periodic events for temporal knowledge graph reasoning". In: *Proceedings of The 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2023, pp. 1559–1568.
- [3] Kaixuan Yao et al. "Multi-view graph convolutional networks with attention mechanism". In: *Artificial Intelligence* 307 (2022), p. 103708.
- [4] Yifan Feng et al. "Hypergraph neural networks". In: *Proceedings of the AAAI Conference on Artificial Intelligence*. 2019, pp. 3558–3565.
- [5] Hongyi Zhang et al. "mixup: Beyond Empirical Risk Minimization". In: *International Conference on Learning Representations*. 2018.

- [6] Lirong Wu et al. "Graphmixup: Improving class-imbalanced node classification by reinforcement mixup and self-supervised context prediction". In: *Joint European Conference on Machine Larning and Knowledge Discovery in Databases*. Springer. 2022, pp. 519–535.
- [7] Sangdoo Yun et al. "Cutmix: Regularization strategy to train strong classifiers with localizable features". In: *Proceedings of The IEEE/CVF International Conference on Computer Vision*. 2019, pp. 6023–6032.
- [8] Shashanka Venkataramanan et al. "Alignmixup: Improving representations by interpolating aligned features". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022, pp. 19174–19183.
- [9] Vikas Verma et al. "Manifold mixup: Better representations by interpolating hidden states". In: *International Conference on Machine Learning*. 2019, pp. 6438–6447.
- [10] Minui Hong, Jinwoo Choi, and Gunhee Kim. "Stylemix: Separating content and style for enhanced data augmentation". In: *Proceedings of The IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 14862–14870.
- [11] Hyeong Kyu Choi, Joonmyung Choi, and Hyunwoo J Kim. "Tokenmixup: Efficient attention-guided token-level data augmentation for transformers". In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 14224–14235.
- [12] JangHyun Kim et al. "Co-Mixup: Saliency guided joint mixup with supermodular diversity". In: *International Conference on Learning Representations*.
- [13] Yuning You et al. "Graph contrastive learning with augmentations". In: *Advances in neural information processing systems* 33 (2020), pp. 5812–5823.
- [14] Xiaotian Han et al. "G-mixup: Graph data augmentation for graph classification". In: *International Conference on Machine Learning*. PMLR. 2022, pp. 8230–8248.
- [15] Vikas Verma et al. "Graphmix: Improved training of gnns for semi-supervised learning". In: Proceedings of The AAAI Conference on Artificial Intelligence. Vol. 35. 11. 2021, pp. 10024– 10032.
- [16] Kaize Ding et al. "Be more with less: Hypergraph attention networks for inductive text classification". In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2020, pp. 4927–4936.
- [17] Song Bai, Feihu Zhang, and Philip HS Torr. "Hypergraph convolution and hypergraph attention". In: *Pattern Recognition* 110 (2021), p. 107637.
- [18] Naganand Yadati et al. "Hypergen: A new method for training graph convolutional networks on hypergraphs". In: *Advances in Neural Information Processing Systems*. 2019.
- [19] Yihe Dong, Will Sawin, and Yoshua Bengio. "HNHN: Hypergraph Networks with Hyperedge Neurons". In: *ICML Graph Representation Learning and Beyond Workshop* (2020). URL: https://arxiv.org/abs/2006.12278.
- [20] Zhenyu Ye, Guangcong Liu, and Zhendong Chen. "Contrastive learning hypergraph neural network based on multimodality replacing data augmentation". In: 2025 IEEE 17th International Conference on Computer Research and Development (ICCRD). IEEE. 2025, pp. 11–16.
- [21] Jian Wang et al. "Hypergraph collaborative filtering with adaptive augmentation of graph data for recommendation". In: *IEEE Transactions on Knowledge and Data Engineering* (2025).
- [22] Olivier Chapelle et al. "Vicinal risk minimization". In: Advances in Neural Information Processing Systems. 2000.
- [23] Prithviraj Sen et al. "Collective classification in network data". In: *AI magazine* 29.3 (2008), pp. 93–93.
- [24] Zhirong Wu et al. "3d shapenets: A deep representation for volumetric shapes". In: *Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition*. 2015, pp. 1912–1920.
- [25] Jiaxin Li, Ben M Chen, and Gim Hee Lee. "So-net: Self-organizing network for point cloud analysis". In: *Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition*. 2018, pp. 9397–9406.
- [26] Hang Su et al. "Multi-view convolutional neural networks for 3d shape recognition". In: *Proceedings of The IEEE International Conference on Computer Vision*. 2015, pp. 945–953.
- [27] Yifan Feng et al. "Gvcnn: Group-view convolutional neural networks for 3d shape recognition". In: *Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition*. 2018, pp. 264–272.

- [28] Petar Velickovic et al. "Graph attention networks". In: *stat* 1050.20 (2017), pp. 10–48550.
- [29] Will Hamilton, Zhitao Ying, and Jure Leskovec. "Inductive representation learning on large graphs". In: *Advances in neural information processing systems* 30 (2017).
- [30] Christopher Morris et al. "Weisfeiler and leman go neural: Higher-order graph neural networks". In: *Proceedings of The AAAI Conference on Artificial Intelligence*. Vol. 33. 01. 2019, pp. 4602–4609.
- [31] Yiwei Wang et al. "Mixup for node and graph classification". In: *Proceedings of the Web Conference 2021*. 2021, pp. 3663–3674.
- [32] Thomas N Kipf and Max Welling. "Semi-Supervised Classification with Graph Convolutional Networks". In: *International Conference on Learning Representations*. 2017.
- [33] Yue Gao et al. "HGNN+: General hypergraph neural networks". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45.3 (2022), pp. 3181–3199.

# **NeurIPS Paper Checklist**

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Yes, the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Yes, we have discussed the limitations of the work in the Conclusion.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

# 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: For each theoretical result, we have provided the full set of assumptions and a complete (and correct) proof.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We have fully disclosed all the information needed to reproduce the main experimental results of the paper

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
  well by the reviewers: Making the paper reproducible is important, regardless of
  whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

# 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We have provided open access to the necessary data and code.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

#### 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We gave pecified all the training and test details.

# Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
  material.

# 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We have reported the compared details of the experiments.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

#### 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We have provided sufficient information on the computer resources.

# Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Yes, the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics.

# Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
  deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

# 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: There is no societal impact of the work performed.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

# 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

# A Detailed Proofs

# A.1 Proof of Theorem 1 (Hypergraph Regularization Decomposition)

To analyze the regularization effect of HyperMixup, we begin by considering the mixed feature for a node pair (i, j), defined as:

$$\tilde{x}_{ij} = \lambda x_i + (1 - \lambda)x_j + \gamma(\lambda x_{e_i} + (1 - \lambda)x_{e_i}),\tag{16}$$

where  $x_{e_i} = D_{e_i}^{-1} \sum_{v \in e_i} x_v$  represents the hyperedge feature. Expanding the loss function around the original feature  $x_i$  through second-order Taylor series yields:

$$l(\theta, \tilde{z}_{ij}) = l(\theta, z_i) + \nabla_x l_i^{\top} (\tilde{x}_{ij} - x_i) + \frac{\partial l_i}{\partial y} (\tilde{y}_{ij} - y_i)$$

$$+ \frac{1}{2} (\tilde{x}_{ij} - x_i)^{\top} \nabla_x^2 l_i (\tilde{x}_{ij} - x_i) + o(\|\tilde{x}_{ij} - x_i\|^2).$$
(17)

Taking expectation over the Beta-distributed mixing coefficient  $\lambda$  and all node pairs (i, j), we decompose the expectation into three components. The linear term involves the gradient alignment between node features and hyperedge deviations:

$$\mathbb{E}\left[\nabla_{x}l_{i}^{\top}(\tilde{x}-x_{i})\right] = \mathbb{E}_{\lambda}[1-\lambda]\nabla_{x}l_{i}^{\top}\left(\mathbb{E}_{j}[x_{j}-x_{i}]+\gamma\mathbb{E}_{j}[x_{e_{j}}-x_{e_{i}}]\right)$$

$$= \mathbb{E}_{\lambda}[1-\lambda]\nabla_{x}l_{i}^{\top}\left(\mu_{v}-x_{i}+\gamma(\mu_{e}-x_{e_{i}})\right),$$
(18)

where  $\mu_v$  and  $\mu_e$  denote the global node and hyperedge feature means, respectively.

The quadratic term captures curvature regularization through the hypergraph covariance structure:

$$\frac{1}{2}\mathbb{E}\left[\left(\tilde{x} - x_i\right)^{\top} \nabla_x^2 l_i (\tilde{x} - x_i)\right] = \frac{\mathbb{E}_{\lambda}\left[\left(1 - \lambda\right)^2\right]}{2} \operatorname{Tr}\left(\nabla_x^2 l_i \cdot \mathbb{E}_{j}\left[\left(\Delta x_{ij} + \gamma \Delta x_{e_{ij}}\right) (\Delta x_{ij} + \gamma \Delta x_{e_{ij}})^{\top}\right]\right) \\
= \frac{\mathbb{E}_{\lambda}\left[\left(1 - \lambda\right)^2\right]}{2} \left(h_i'' \nabla f_i^{\top} (\Sigma_v + \gamma^2 \Sigma_e) \nabla f_i + \gamma (h_i' - y_i) \operatorname{Tr}(\nabla^2 f_i \Sigma_{ve})\right), \tag{19}$$

where  $\Delta x_{ij} = x_j - x_i$ ,  $\Delta x_{e_{ij}} = x_{e_j} - x_{e_i}$ , and  $\Sigma_v$ ,  $\Sigma_e$ ,  $\Sigma_{ve}$  represent node covariance, hyperedge covariance, and their cross-covariance matrices.

Under the uniform hyperedge sampling assumption  $\mathbb{E}_j[x_{e_j}] = \mu_e$ , the cross-covariance  $\Sigma_{ve}$  vanishes, simplifying the expression to the stated regularization terms  $\mathcal{R}_1$ ,  $\mathcal{R}_2$ , and  $\mathcal{R}_3$ . This decomposition explicitly reveals how HyperMixup introduces hypergraph-aware regularization through 1) nodehyperedge gradient alignment, 2) hyperedge covariance-driven curvature penalization, and 3) higher-order interactions between loss Hessian and hyperedge structure.

#### A.2 Proof of Theorem 2 (Hypergraph Robustness Bound)

Consider adversarial perturbations  $\delta = \delta_v + \gamma \delta_e$  affecting both node features and hyperedge propagations, bounded by  $\|\delta_v\|_2 \le \epsilon_v \sqrt{d}$  and  $\|\delta_e\|_2 \le \epsilon_e \sqrt{d}$ . Expanding the loss difference for perturbed features  $x_i' = x_i + \delta$  gives:

$$l(\theta, x_i') - l(\theta, x_i) = \nabla_x l_i^{\top} (\delta_v + \gamma \delta_e) + \frac{1}{2} (\delta_v + \gamma \delta_e)^{\top} \nabla_x^2 l_i (\delta_v + \gamma \delta_e) + o(\|\delta\|^2).$$
 (20)

Maximizing over admissible perturbations reveals the worst-case loss increase:

$$\max_{\delta} l(\theta, x_i') \le l(\theta, x_i) + \epsilon_v \sqrt{d} \|\nabla_x l_i\| + \gamma \epsilon_e \sqrt{d} \|\nabla_x l_i\| + \frac{d}{2} (\epsilon_v^2 + \gamma^2 \epsilon_e^2) \lambda_{\max}(\nabla_x^2 l_i) + o(d).$$
(21)

Relating this to the HyperMixup regularization terms derived in Theorem 1, we observe that  $\mathcal{R}_1$  controls the linear gradient norms through  $\mathbb{E}_{\lambda}[1-\lambda]\|\nabla f_i\|$ , while  $\mathcal{R}_2$  and  $\mathcal{R}_3$  constrain the Hessian spectral norm  $\lambda_{\max}(\nabla_x^2 l_i)$ . The effective perturbation radius  $\epsilon_{\min}$  emerges as a weighted combination of node and hyperedge attack strengths, scaled by the alignment factor  $R=\min_i|\cos(\nabla f_i,x_e-x_i)|$  that quantifies consistency between node gradients and hyperedge structure.

This analysis rigorously establishes that HyperMixup training minimizes an upper bound of the adversarial loss, with the hyperparameter  $\gamma$  dynamically balancing robustness between node-level and hyperedge-level attacks. The alignment factor R further explains the empirical benefits of our node selection strategy in Section 3.3.1, which explicitly maximizes gradient-hyperedge consistency.

# A.3 Proof of Theorem 3 (Hypergraph Generalization)

The Rademacher complexity analysis begins with the spectral decomposition of the hypergraph Laplacian  $L = U\Lambda U^{\top}$ , where U contains the spectral basis vectors. Expressing the model function in this basis:

$$f_{\theta}(x_i) = \sum_{k=1}^{K} \theta_k u_k(i), \tag{22}$$

we bound the complexity through spectral energy concentration:

$$\operatorname{Rad}_{n}(\mathcal{F}_{\mathcal{G}}) = \mathbb{E}_{\xi} \left[ \sup_{\|\theta\|_{\mathcal{G}} \leq B} \frac{1}{n} \sum_{i=1}^{n} \xi_{i} \sum_{k=1}^{K} \theta_{k} u_{k}(i) \right]$$

$$\leq \frac{B}{n} \mathbb{E}_{\xi} \left[ \sqrt{\sum_{k=1}^{K} \left( \sum_{i=1}^{n} \xi_{i} u_{k}(i) \right)^{2}} \right]$$

$$\leq \frac{B}{n} \sqrt{n} \sum_{k=1}^{K} \|u_{k}\|_{2}^{2}$$

$$= B\sqrt{\frac{K}{n}}. \tag{23}$$

The effective dimension K is constrained by hypergraph spectral properties:

$$K \le \rho(L) \left( r + \gamma^2 \| \Sigma_{ve} \|_F^2 \right), \tag{24}$$

where  $\rho(L)$  denotes the spectral radius encoding hypergraph connectivity,  $r = \operatorname{rank}(\Sigma_v)$  reflects node feature dimensionality, and  $\|\Sigma_{ve}\|_F$  quantifies node-hyperedge feature alignment. Substituting this into the complexity bound yields the final result:

$$\operatorname{Rad}_{n}(\mathcal{F}_{\mathcal{G}}) \leq \sqrt{\frac{C\rho(L)(r + \gamma^{2} \|\Sigma_{ve}\|_{F}^{2})}{n}}.$$
(25)

This bound reveals the generalization benefits of HyperMixup: 1) The spectral radius  $\rho(L)$  encourages adaptation to hypergraph density through the  $\mathcal{R}_2$  regularization; 2) The  $\gamma^2 \|\Sigma_{ve}\|_F^2$  term formalizes the advantage of hyperedge mixing when node features align with hyperedge structure; 3) Low-rank node covariance r (typical in citation networks) naturally reduces model complexity. These theoretical insights align with the empirical observations in Table 2, particularly the superior performance on Cora compared to ModelNet40.