

Graph Your Own Prompt

Xi Ding^{1,3}, Lei Wang^{1,2}, Piotr Koniusz^{1,2,3,4}, Yongsheng Gao¹

¹Griffith University

²Data61/CSIRO

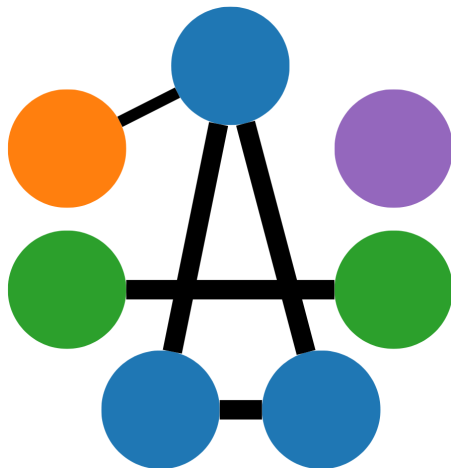
³Australian National University

⁴University of New South Wales



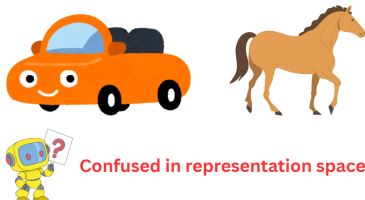
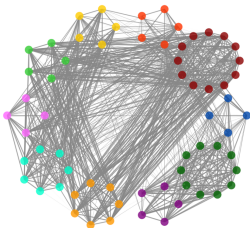
Outline

- 1 Motivation
- 2 Method
- 3 Experiments
- 4 Discussion
- 5 Conclusion



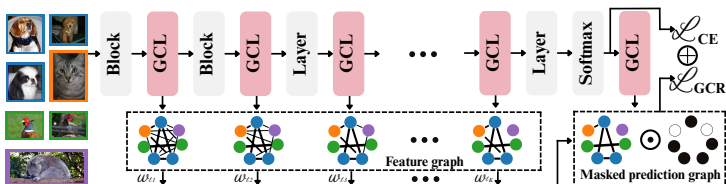
Motivation: Why Structure the Feature Space?

- Deep neural networks achieve strong accuracy but their internal features are often:
 - noisy and entangled,
 - misaligned with semantic class boundaries,
 - hard to interpret and fragile under distribution shifts.
- In many trained models, samples from different classes can remain close in feature space, while same-class samples may not cluster well.
- This harms generalization, robustness and interpretability, even when the final classifier performs well.



Our Idea in One Slide

- Build relational graphs directly from:
 - intermediate features and
 - class predictions within each mini-batch.
- Use the prediction graph as a semantic reference.
- Regularize the network so that feature graphs at multiple layers align with this semantic prediction graph.
- This acts as an internal “prompt”:
 - the model uses its own prediction structure to refine its feature geometry,
 - without changing the backbone architecture or adding trainable parameters.



Overview of Graph Consistency Regularization (GCR)

- **Goal:** enforce consistency between
 - batch-level feature similarity graphs and
 - a class-aware prediction similarity graph.
- **Notation (per batch):**
 - feature matrix at layer l : $X^{(l)} \in \mathbb{R}^{n \times d}$, with rows $x_i^{(l)}$,
 - logits: $Z \in \mathbb{R}^{n \times C}$,
 - ground-truth labels: y_1, \dots, y_n .
- **Key components:**
 - Graph Consistency Layers (GCLs) inserted after chosen blocks or layers,
 - a masked prediction relational graph P built from softmax outputs and labels,
 - a graph alignment loss aggregated across layers with adaptive weights.
- **Result:** semantically structured features, stronger intra-class cohesion and reduced noisy inter-class affinities.

Graph Consistency Layer: Feature Graph

- At a chosen layer l , collect the batch features as a matrix:

$$\mathbf{X}^{(l)} = \begin{bmatrix} (\mathbf{x}_1^{(l)})^\top \\ \vdots \\ (\mathbf{x}_n^{(l)})^\top \end{bmatrix} \in \mathbb{R}^{n \times d},$$

where each row is a sample flattened into a feature vector.

- Build a feature similarity graph $\mathbf{F}^{(l)} \in \mathbb{R}^{n \times n}$:
 - nodes correspond to samples in the batch,
 - edges encode pairwise similarity between feature vectors.
- We use cosine similarity with non-negative values:

$$F_{ij}^{(l)} = \text{ReLU}(\cos(\mathbf{x}_i^{(l)}, \mathbf{x}_j^{(l)})), \quad i, j = 1, \dots, n. \quad (1)$$

- This graph captures how the model currently organizes samples in feature space at that layer.

Graph Consistency Layer: Masked Prediction Graph

- From the prediction logits $Z = [z_1^\top, \dots, z_n^\top]^\top$ of the same batch:
 - apply softmax to obtain class probability vectors $p_i = \text{softmax}(z_i)$,
 - compute pairwise cosine similarity between prediction vectors:

$$S_{ij} = \text{ReLU}(\cos(p_i, p_j)). \quad (2)$$

- To focus on reliable semantic relations, we build a binary mask $M \in \{0, 1\}^{n \times n}$:

$$M_{ij} = \begin{cases} 1, & \text{if } y_i = y_j, \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

- The masked prediction graph $P \in \mathbb{R}^{n \times n}$ is then

$$P_{ij} = M_{ij} \odot S_{ij}, \quad (4)$$

where \odot denotes elementwise multiplication.

- This graph keeps intra-class semantic structure and discards possibly misleading inter-class similarities.

Layer-wise Graph Alignment

- For each GCL at layer l , we align:
 - the feature graph $F^{(l)}$ at that layer,
 - with the fixed masked prediction graph P for the batch.
- We:
 - keep only the strictly upper triangular part (undirected graph, no self-loops):

$\text{triu}(A)$ = upper triangular part of A , without diagonal,

- measure the discrepancy via the squared Frobenius norm.
- The layer-wise **graph consistency loss** is

$$\mathcal{L}_{\text{GCR}}^{(l)} = \|\text{triu}(F^{(l)}) - \text{triu}(P)\|_F^2. \quad (5)$$

- Small loss means feature relationships follow the semantic relationships implied by predictions; large loss means mismatch between geometry and semantics.

Aggregating Across Layers

- GCL can be placed at multiple depths: early, middle, late, or all stages.
- For a set of layers $\{1, \dots, K\}$, compute a graph consistency loss at each layer and combine them:

$$\mathcal{L}_{\text{GCR}} = \sum_{l=1}^K w_l \|\text{triu}(\mathbf{F}^{(l)}) - \text{triu}(\mathbf{P})\|_F^2, \quad (6)$$

where $w_l \geq 0$ are layer weights.

- **Fixed weighting** examples:

$$w_l \in \left\{ \frac{1}{K}, \frac{l}{K}, \left(\frac{l}{K}\right)^2, \frac{1+\cos(\pi l/K)}{2}, \frac{\arccos(1-2l/K)}{\pi} \right\}.$$

- **Adaptive weighting** (graph-discrepancy-based):

$$w_l = \frac{\exp(-\|\text{triu}(\mathbf{F}^{(l)}) - \text{triu}(\mathbf{P})\|_F^2)}{\sum_{j=1}^K \exp(-\|\text{triu}(\mathbf{F}^{(j)}) - \text{triu}(\mathbf{P})\|_F^2)}. \quad (7)$$

Training Objective

- The standard supervised objective is cross-entropy on the predictions:

$$\mathcal{L}_{\text{CE}}.$$

- Our total loss adds the GCR term:

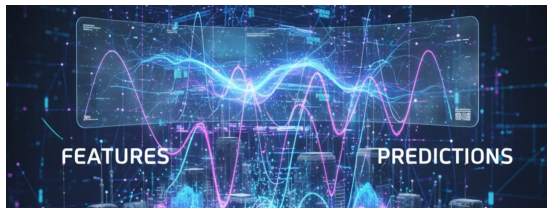
$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CE}} + \lambda \mathcal{L}_{\text{GCR}}, \quad (8)$$

where $\lambda > 0$ controls the strength of the regularization.

- Interpretation:
 - \mathcal{L}_{CE} fits labels at the prediction level,
 - \mathcal{L}_{GCR} aligns feature-space relations with class-aware prediction structure.
- Implementation is **portable**:
 - no new learnable parameters,
 - only involves matrix operations over the current batch,
 - can be applied to many backbones without redesigning the architecture.

GCR as Self-Prompting

- Traditional prompting:
 - injects external tokens or instructions into the model.
- GCR can be viewed as **internal prompting**:
 - the model uses its own masked prediction graph P to generate a structural signal,
 - this signal shapes feature relationships $F^{(l)}$ across layers via the loss in (5)–(8).
- Characteristics of this internal prompt:
 - purely internal (no external tokens),
 - structural (operates on pairwise relations, not single samples),
 - semantic (class-aware masking makes the graph class-consistent).



Echoes in the Model: When Features Reflect Predictions

- **Datasets:**

- Kaggle Cats vs. Dogs,
- CIFAR-10 and CIFAR-100,
- Tiny ImageNet,
- ImageNet-1K.

- **Architectures:**

- lightweight CNNs: MobileNet, ShuffleNet, SqueezeNet, GoogLeNet,
 - deeper CNNs: ResNet, ResNeXt, DenseNet, stochastic and squeeze-and-excitation variants,
 - vision transformers: ViT, Swin, MobileViT, CEiT, iFormer, ViG,
 - masked autoencoders.
- GCLs are inserted at early, middle, late stages, or combinations of them.

Quantitative Results on CIFAR-10/100

	MAE	MNet	SN	SQNet	GLNet	Rx-50	Rx-101	R34	R50	R101	D121	Mean
Baseline	88.95 \pm 0.33	90.23 \pm 0.25	91.21 \pm 0.28	92.30 \pm 0.25	94.10 \pm 0.26	94.57 \pm 0.29	95.12 \pm 0.30	94.83 \pm 0.25	95.03 \pm 0.28	95.22 \pm 0.31	95.01 \pm 0.27	93.32 \pm 2.26
Early GCL	89.42 \pm 0.25	91.17 \pm 0.22	92.33 \pm 0.33	92.59 \pm 0.21	94.89 \pm 0.23	95.48 \pm 0.22	95.63 \pm 0.25	95.55 \pm 0.18	95.57 \pm 0.23	95.39 \pm 0.26	95.81 \pm 0.17	93.98 \pm 2.22
Mid GCL	89.77 \pm 0.22	91.15 \pm 0.18	92.58 \pm 0.19	92.40 \pm 0.20	94.82 \pm 0.21	95.47 \pm 0.19	95.39 \pm 0.24	95.69 \pm 0.23	95.61 \pm 0.20	95.75 \pm 0.17	95.51 \pm 0.22	94.01 \pm 2.15
Late GCL	89.70 \pm 0.29	91.40 \pm 0.19	92.36 \pm 0.21	92.80 \pm 0.19	94.88 \pm 0.19	95.35 \pm 0.28	95.71 \pm 0.26	95.69 \pm 0.19	95.66 \pm 0.17	95.51 \pm 0.24	95.72 \pm 0.22	94.07 \pm 2.14
Early+Mid	89.52 \pm 0.19	90.77 \pm 0.26	92.56 \pm 0.21	92.27 \pm 0.25	94.79 \pm 0.18	95.33 \pm 0.27	95.55 \pm 0.23	95.46 \pm 0.20	95.51 \pm 0.21	95.37 \pm 0.19	95.64 \pm 0.20	93.89 \pm 2.22
Mid+Late	89.59 \pm 0.28	91.23 \pm 0.20	92.79 \pm 0.20	92.86 \pm 0.23	94.61 \pm 0.22	95.51 \pm 0.19	95.38 \pm 0.27	95.45 \pm 0.18	95.33 \pm 0.26	95.52 \pm 0.14	95.70 \pm 0.19	94.00 \pm 2.09
Early+Late	89.64 \pm 0.25	91.03 \pm 0.24	92.30 \pm 0.28	92.70 \pm 0.23	94.69 \pm 0.20	95.40 \pm 0.20	95.35 \pm 0.23	95.66 \pm 0.21	95.31 \pm 0.25	95.49 \pm 0.16	95.53 \pm 0.22	93.92 \pm 2.14
Full GCL	89.55 \pm 0.23	90.99 \pm 0.18	92.48 \pm 0.19	92.65 \pm 0.20	94.57 \pm 0.21	95.50 \pm 0.19	95.34 \pm 0.20	95.48 \pm 0.17	95.62 \pm 0.18	95.38 \pm 0.21	95.51 \pm 0.20	93.92 \pm 2.15

Accuracy (%) on CIFAR-10 across models.

	MAE	MNet	SN	SQNet	Rx-50	Rx-101	R34	R50	D121	Mean
Baseline	64.29 \pm 0.34	65.95 \pm 0.25	70.11 \pm 0.30	69.43 \pm 0.27	77.75 \pm 0.29	77.83 \pm 0.30	76.82 \pm 0.28	77.31 \pm 0.29	77.09 \pm 0.27	72.95 \pm 5.50
Early GCL	65.05 \pm 0.29	67.45 \pm 0.21	71.96 \pm 0.27	70.90 \pm 0.20	79.18 \pm 0.22	79.69 \pm 0.27	77.90 \pm 0.22	79.37 \pm 0.25	79.41 \pm 0.22	74.55 \pm 5.78
Mid GCL	64.99 \pm 0.30	67.88 \pm 0.21	71.89 \pm 0.24	70.21 \pm 0.25	79.07 \pm 0.19	79.28 \pm 0.26	77.83 \pm 0.20	78.90 \pm 0.24	79.26 \pm 0.21	74.37 \pm 5.66
Late GCL	65.54 \pm 0.27	68.32 \pm 0.20	71.42 \pm 0.24	70.55 \pm 0.22	79.54 \pm 0.20	79.83 \pm 0.21	78.31 \pm 0.20	79.42 \pm 0.21	79.69 \pm 0.23	74.74 \pm 5.73
Early+Mid	65.23 \pm 0.31	67.62 \pm 0.24	71.50 \pm 0.28	70.47 \pm 0.19	78.90 \pm 0.18	79.25 \pm 0.20	77.41 \pm 0.19	78.58 \pm 0.24	79.22 \pm 0.20	74.28 \pm 5.56
Mid+Late	65.27 \pm 0.28	68.33 \pm 0.19	71.63 \pm 0.28	70.30 \pm 0.22	78.91 \pm 0.17	79.57 \pm 0.21	77.30 \pm 0.20	78.85 \pm 0.22	79.54 \pm 0.24	74.41 \pm 5.55
Early+Late	65.22 \pm 0.21	67.25 \pm 0.21	71.55 \pm 0.27	71.03 \pm 0.24	79.03 \pm 0.20	79.41 \pm 0.22	78.19 \pm 0.23	78.70 \pm 0.23	79.45 \pm 0.22	74.43 \pm 5.69
Full GCL	65.38 \pm 0.22	68.22 \pm 0.19	71.30 \pm 0.24	70.77 \pm 0.20	79.01 \pm 0.19	79.29 \pm 0.21	77.79 \pm 0.20	78.71 \pm 0.22	79.27 \pm 0.19	74.42 \pm 5.49

Accuracy (%) on CIFAR-100 across models.

Quantitative Results on Tiny ImageNet/ImageNet-1K

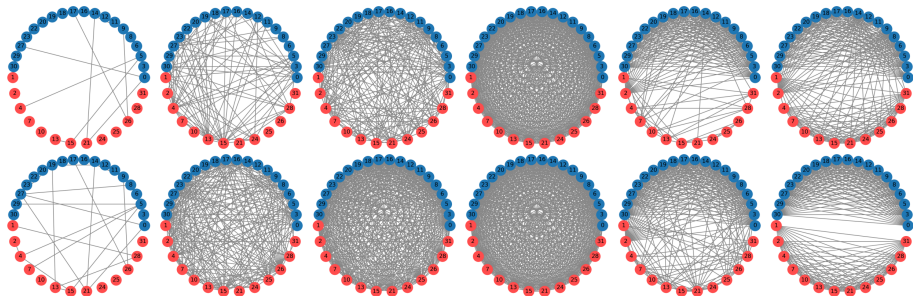
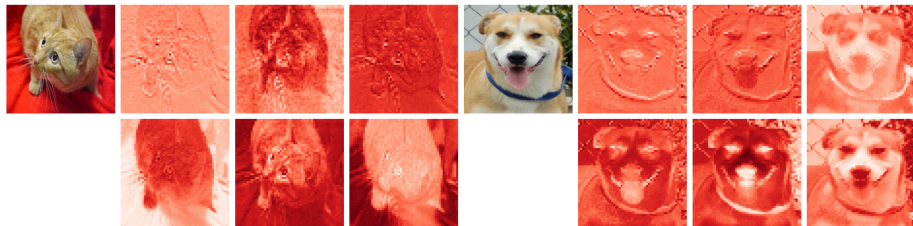
	ViT _{/32}	ViT _{/16}	CeT	MViT _{XCS}	MViT _{XS}	MViT	Swin	MNet	R18SD	SER18	R34	Mean
Baseline	37.79 \pm 0.35	40.05 \pm 0.33	49.95 \pm 0.29	49.28 \pm 0.29	51.58 \pm 0.27	52.68 \pm 0.27	54.27 \pm 0.25	57.81 \pm 0.25	63.49 \pm 0.26	65.65 \pm 0.24	67.51 \pm 0.25	53.64 \pm 0.62
Early GCL	39.02 \pm 0.29	40.98 \pm 0.19	51.22 \pm 0.20	50.11 \pm 0.28	51.33 \pm 0.26	53.91 \pm 0.22	54.88 \pm 0.25	57.93 \pm 0.21	63.81 \pm 0.19	66.52 \pm 0.22	67.79 \pm 0.19	54.32 \pm 0.39
Mid GCL	38.61 \pm 0.23	40.95 \pm 0.19	50.30 \pm 0.19	49.92 \pm 0.26	51.43 \pm 0.22	53.88 \pm 0.20	55.23 \pm 0.24	57.63 \pm 0.20	64.03 \pm 0.22	65.66 \pm 0.23	67.62 \pm 0.20	54.11 \pm 0.38
Late GCL	37.98 \pm 0.28	40.35 \pm 0.25	<u>50.82</u> \pm 0.20	<u>49.77</u> \pm 0.21	51.99 \pm 0.23	54.10 \pm 0.19	<u>55.47</u> \pm 0.21	57.87 \pm 0.23	63.79 \pm 0.19	65.85 \pm 0.25	67.61 \pm 0.19	54.18 \pm 0.56
Early+Mid	39.08 \pm 0.25	41.26 \pm 0.18	50.25 \pm 0.25	49.73 \pm 0.22	51.57 \pm 0.19	53.91 \pm 0.23	54.95 \pm 0.19	57.49 \pm 0.19	<u>64.18</u> \pm 0.20	65.86 \pm 0.24	67.74 \pm 0.23	54.18 \pm 0.32
Mid+Late	38.44 \pm 0.18	40.52 \pm 0.28	50.09 \pm 0.25	50.55 \pm 0.18	51.48 \pm 0.21	53.90 \pm 0.20	55.62 \pm 0.23	57.65 \pm 0.21	64.29 \pm 0.17	65.95 \pm 0.19	67.58 \pm 0.21	54.19 \pm 0.52
Early+Late	38.34 \pm 0.23	40.71 \pm 0.21	50.70 \pm 0.25	50.23 \pm 0.20	51.36 \pm 0.18	53.57 \pm 0.21	54.89 \pm 0.21	57.93 \pm 0.19	63.88 \pm 0.19	65.83 \pm 0.17	<u>67.75</u> \pm 0.25	<u>54.11</u> \pm 0.47
Full GCL	38.38 \pm 0.22	40.80 \pm 0.18	49.92 \pm 0.28	50.16 \pm 0.17	<u>51.87</u> \pm 0.19	<u>54.01</u> \pm 0.19	54.87 \pm 0.19	57.64 \pm 0.20	<u>64.10</u> \pm 0.19	<u>66.01</u> \pm 0.15	67.66 \pm 0.18	54.13 \pm 0.49

Accuracy (%) on Tiny ImageNet across models.

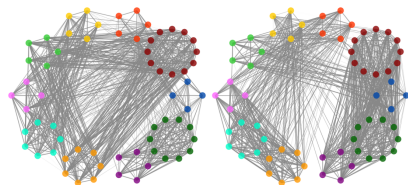
Method	iFormer-S	iFormer-B	ViT-B/16	ViG-B
Baseline	83.4 \pm 0.40	84.6 \pm 0.45	74.3 \pm 0.51	82.3 \pm 0.42
Early GCL	83.8 \pm 0.31	85.0 \pm 0.40	74.7 \pm 0.44	82.8 \pm 0.35
Mid GCL	83.8 \pm 0.39	85.5 \pm 0.33	75.2 \pm 0.36	83.0 \pm 0.34
Late GCL	<u>84.5</u> \pm 0.29	86.1 \pm 0.30	75.8 \pm 0.33	84.0 \pm 0.30
Early + Mid	84.3 \pm 0.33	85.9 \pm 0.38	<u>75.6</u> \pm 0.41	83.7 \pm 0.33
Mid + Late	84.8 \pm 0.28	<u>85.9</u> \pm 0.37	<u>75.6</u> \pm 0.34	<u>83.9</u> \pm 0.30
Early + Late	<u>84.5</u> \pm 0.30	85.2 \pm 0.28	74.9 \pm 0.33	83.5 \pm 0.29
Full GCL	84.3 \pm 0.29	85.8 \pm 0.26	75.5 \pm 0.30	83.6 \pm 0.27

Accuracy (%) on ImageNet-1K across models.

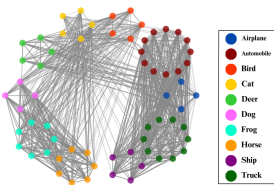
Qualitative Results on Kaggle Cats vs. Dogs



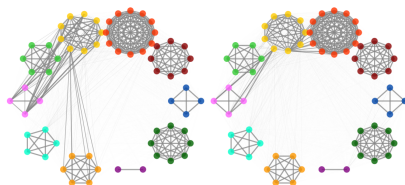
Qualitative Results on CIFAR-10



(a) DenseNet-121

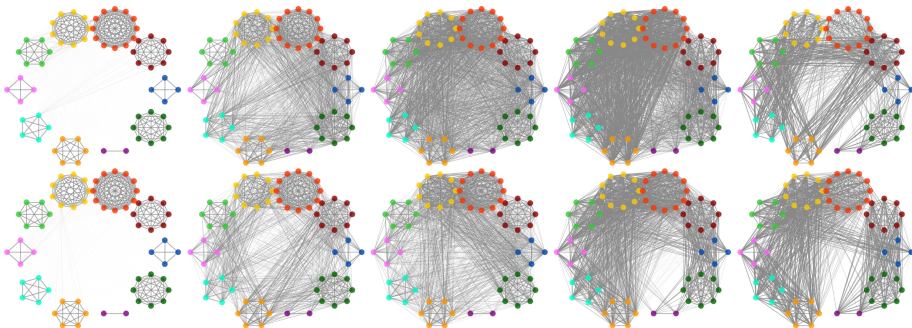


(b) With our GCLs



(c) MobileNet

(d) With our GCLs



(a) ShuffleNet

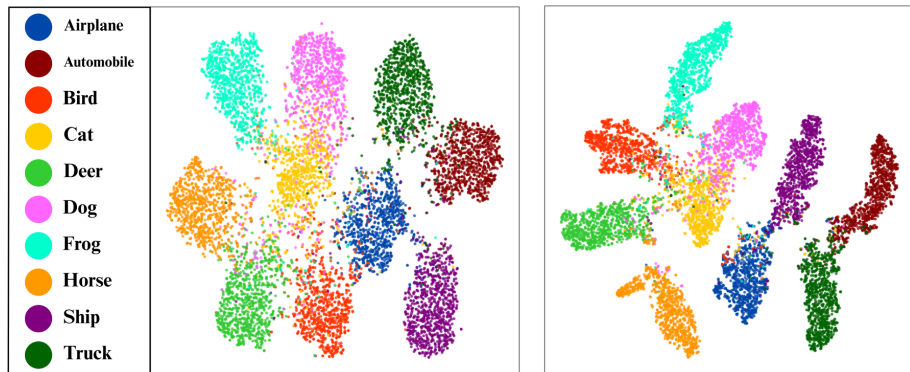
(b) GoogLeNet

(c) ResNet-50

(d) ResNet-101

(e) DenseNet-121

Qualitative Results on CIFAR-10



Left is the baseline (ShuffleNet), and the right is GCL-augmented. GCL-augmented yields more **compact intra-class clusters** and **better inter-class separation**.

Where to Place & How to Weight GCLs?

Linear	0.55	0.45	0.48	0.31	0.34	0.36	0.27
Sqrt	0.44	0.41	0.42	0.21	0.33	0.26	0.19
Squared	0.55	0.46	0.49	0.51	0.49	0.47	0.40
Equal	0.49	0.58	0.56	0.30	0.47	0.45	0.45
Arccos	0.51	0.53	0.53	0.34	0.40	0.40	0.24
Cosine	0.39	0.44	0.57	0.23	0.45	0.39	0.19
Adaptive	0.48	0.56	0.63	0.45	0.61	0.53	0.51
	E	M	L	E+M	M+L	E+L	Full

CIFAR-10

Linear	1.40	1.20	1.30	0.80	0.96	0.95	0.77
Sqrt	1.15	1.02	1.08	0.66	0.75	0.79	0.40
Squared	1.47	1.16	1.14	1.12	0.90	1.14	0.97
Equal	1.10	1.11	1.24	1.14	1.37	1.11	1.26
Arccos	1.38	1.17	0.98	0.85	0.76	0.92	0.53
Cosine	1.17	1.23	1.52	0.60	1.14	1.22	0.67
Adaptive	1.45	1.30	1.57	1.27	1.23	1.42	1.30
	E	M	L	E+M	M+L	E+L	Full

CIFAR-100

● Placement:

- Late GCLs yield the largest gains.
- Early GCLs help stabilize feature learning and reduce noise propagation.
- Combining late with other stages is often beneficial.
- Full GCL not always best (too much regularization may hurt).

● Weighting:

- Adaptive weighting works best (focuses on misaligned layers).
- Fixed schemes: squared / equal > linear / square root.

What Does GCR Buy Us?

- **Semantic feature geometry:**

- features are organized to respect class-aware relationships in prediction space,
- intra-class cohesion increases, inter-class confusion decreases.

- **Better generalization:**

- graph-based regularization shrinks the effective hypothesis space,
- empirical results show consistent gains across datasets and architectures.

- **Interpretability:**

- relational graphs become easier to interpret,
- feature maps highlight semantically meaningful regions.

GCR as a Bridge Between Paradigms

- **Graph-based learning:**

- GCR builds graphs dynamically within each batch,
- no need for persistent global graphs or specialized graph neural networks.

- **Contrastive-style supervision:**

- graph alignment can be seen as a soft form of contrast,
- encourages similar samples to be close and dissimilar samples to be separated,
- without explicit positive and negative sampling.

- **Prompting and self-conditioning:**

- prediction graphs act as continuous, internal prompts that shape representation learning.

When Is GCR Most Helpful?

- Datasets with many classes or fine-grained distinctions:
 - benefit from stronger control over feature geometry,
 - GCR helps separate visually similar but semantically different classes.
- Architectures with limited capacity:
 - lightweight models gain substantial accuracy from the extra structural supervision.
- Deeper backbones:
 - late layers encode rich semantics, and graph alignment at these layers is especially effective.

Limitations and Practical Considerations

- Requires labels to build the class-aware mask:
 - the current formulation is fully supervised,
 - extensions to semi-supervised or self-supervised settings remain open work.
- Graph construction is batch-based:
 - complexity grows with batch size,
 - in practice, the cost is manageable due to efficient matrix operations.
- Strong reliance on prediction quality:
 - very early in training predictions may be noisy,
 - however, the class-aware mask mitigates the influence of unreliable inter-class similarities.

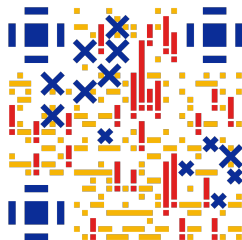
- Introduced **Graph Consistency Regularization (GCR)**:
 - a parameter-free, model-agnostic framework,
 - aligns batch-level feature graphs with a class-aware prediction graph.
- Proposed **Graph Consistency Layers (GCLs)**:
 - lightweight modules that can be inserted at arbitrary depths,
 - enforce multi-layer structural supervision with flexible weighting.
- Demonstrated:
 - improved semantic structure of feature spaces,
 - stronger intra-class cohesion and clearer class boundaries,
 - consistent accuracy gains across many models and datasets.

- Extend GCR beyond image classification:
 - semantic and instance segmentation,
 - retrieval and metric learning tasks,
 - multi-modal settings such as vision-language models.
- Combine GCR with self-supervised and semi-supervised learning:
 - use pseudo-labels or clustering to build masked prediction graphs without full supervision.
- Explore more advanced graph constructions:
 - adaptive sparsification,
 - other similarity measures and masking strategies,
 - alternative ways to incorporate temporal or multi-view structure.

Thank You



Paper



Code

Scan the QR codes to read our paper and code.