

## Mixing Configurations for Downstream Prediction

Anonymous Author(s)

Affiliation  
Address  
email

## Abstract

Humans possess an innate ability to group objects by similarity—a cognitive mechanism that clustering algorithms aim to emulate. Recent advances in community detection have enabled the discovery of *configurations*—valid hierarchical clusterings across multiple resolution scales—without requiring labeled data. In this paper, we formally characterize these configurations and identify similar emergent structures in register tokens within Vision Transformers. Unlike register tokens, configurations exhibit lower redundancy and eliminate the need for ad hoc selection. They can be learned through unsupervised or self-supervised methods, yet their selection or composition remains specific to the downstream task and input. Building on these insights, we introduce GraMixC, a plug-and-play module that extracts configurations, aligns them using our novel Reverse Merge/Split (RMS) technique, and fuses them via attention heads before forwarding them to any downstream predictor. On the DSNI 16S rRNA cultivation-media prediction task, GraMixC improves the  $R^2$  from 0.6 to 0.9 on various methods, setting a new state-of-the-art. We further validate GraMixC across standard tabular benchmarks, where it consistently outperforms single-resolution and static-feature baselines.

17 1 Introduction

Learning general-purpose features that enhance downstream tasks has been a long-standing goal in machine learning. One prominent example is clustering (*i.e.*, community detection) in unsupervised learning, which groups entities into clusters of similar objects while separating dissimilar ones, without using labels. [1, 2, 3]. Interestingly, this paradigm demonstrates remarkable similarities to human-like behaviors. Decades of cognitive science studies show that even infants have the ability to group objects by similarity [4, 5]. In particular, they often organize them at different abstraction levels [6, 7]. Inspired by this, recent advances in community detection have extended clustering to the discovery of *configurations*—hierarchical clusterings that span multiple resolution scales [8]. For example, as illustrated in the lineage diagram of Fig. 1, in the CIFAR10 dataset [9], coarse configurations may separate vehicles from animals, while finer configurations distinguish between birds, cats, and dogs. These multi-resolution representations reveal rich hierarchical structures that could be useful for deep models. However, despite their potential, such configurations have not been widely adopted in deep learning, especially in challenging domains where the underlying data distributions are complex and non-stationary.

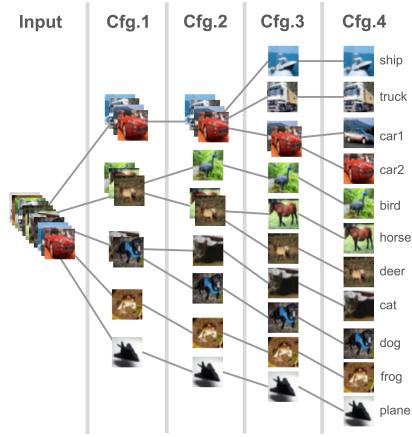


Figure 1: Illustration of CIFAR10 configurations. Each column represents a configuration—**clustering at a specific resolution**.

40 One such domain is 16S ribosomal RNA (rRNA) gene sequencing, a widely used tool in microbiome  
41 studies for identifying and classifying bacteria. Analyzing 16S rRNA data has consistently confronted  
42 significant challenges in downstream prediction tasks within label-scarce environments. Previous  
43 works in 16S rRNA representation learning have demonstrated substantial benefits for bacterial  
44 taxonomic profiling and microbial community analysis [10, 11, 12]. Notably, Johnson et al. [13]  
45 showed that full-length sequencing combined with appropriate clustering of intragenomic sequence  
46 variation can provide more accurate representation of bacterial species in microbiome datasets. These  
47 findings underscore the importance of learning clustered representations without relying on labels.

48 Recent methodologies typically transform clustering results into pseudo-labels to enhance down-  
49 stream prediction performance. For instance, DeepCluster [14] iteratively clusters CNN-extracted  
50 visual features and leverages these cluster assignments to guide network parameter updates. Graph-  
51 based methods such as [15] employ structural clustering to overcome limitations of traditional  
52 contrastive learning approaches that depend on positive and negative sample pairs. Their method  
53 captures structural relationships among nodes in heterogeneous information networks, establishing a  
54 self-supervised pre-training framework that learns robust network representations from unlabeled  
55 data. Nevertheless, aforementioned approaches predominantly focus on a single configuration type,  
56 overlooking the potential benefits of mixing configurations across multiple resolution scales.

57 In this paper, we introduce GraMixC, a plug-and-play module that extracts, aligns and mixes graph-  
58 based configurations for downstream prediction. The main contributions of the paper are as follows:

- 59 • We identify three key characteristics of clustering configurations through systematic exper-  
60 imental analysis, providing a novel perspective on enhancing downstream prediction via  
61 mixing configurations.
- 62 • We propose GraMixC, a plug-and-play module based on mixed configurations. We apply it  
63 to a novel 16S rRNA cultivation-media prediction task, setting a new state-of-the-art.
- 64 • We further conduct extensive experiments on multiple standard tabular benchmarks to  
65 validate GraMixC’s effectiveness, where it consistently outperforms single-resolution and  
66 static-feature baselines.

67 The remainder of this paper is organized as follows. Section 2 analyzes behavioral patterns of  
68 configurations. Section 3 details our proposed GraMixC. Section 4 evaluates GraMixC’s performance  
69 through extensive experiments. Finally, Section 5 concludes the paper. Our data and implementation  
70 is available at <https://anonymous.4open.science/r/project-34CB>.

## 71 2 Preliminary results

72 We first present preliminary experimental results on configurations using CIFAR10. Specifically,  
73 we compare patterns of configurations with those of the learnable “register” tokens in a recent  
74 vision transformer DINoV2-reg [16]. Fig. 2 shows the attention maps from our configurations and  
75 their register tokens. Moreover, Fig. 3 shows qualitative behaviors of our configurations and their  
76 quantitative advantages over registers in terms of feature importance and neighborhood similarity.  
77 From these results, we identify three key properties:

78 **Configurations emerge via unsupervised or self-supervised learning.** We define Near ground truth  
79 (GT) balls as balls selected with the highest clustering scores, marked yellow in Fig. 2a. As shown in  
80 Fig. 2b, the attention map, acquired by feeding configurations as tokens to attention heads for linear  
81 probing, yields high norm regions substantially overlap with GT balls. On another hand, DINoV2-reg  
82 exhibits similar attention map patterns in selected registers (see Fig. 2c), which might be related to  
83 registers activating different areas in Fig. 2d, similar to slot attention [16, 17, 18, 19]. Thus, based on  
84 the similar attention map behavior, register token can be considered as a latent configuration.

85 **Configurations are selected and mixed based on input and task.** *Configuration selection and*  
86 *mixing* refers to learning which resolution scales to focus on for a given downstream task. We  
87 visualize this via attention maps over configuration tokens, where high-norm regions indicate the  
88 selected scales. In Fig. 2b, attention norms vary across rows, showing that each input sample triggers  
89 different resolution scales. Without any change to the configurations, we merge the original labels  
90 into coarser classes (Fig. 3a) and plot the new attention map (Fig. 3b). The attention shifts to align  
91 with the coarser GT, whereas DINoV2-reg register tokens remain unchanged unless re-trained. These  
92 observations confirm that configuration selection and mixing are input- and task-dependent.

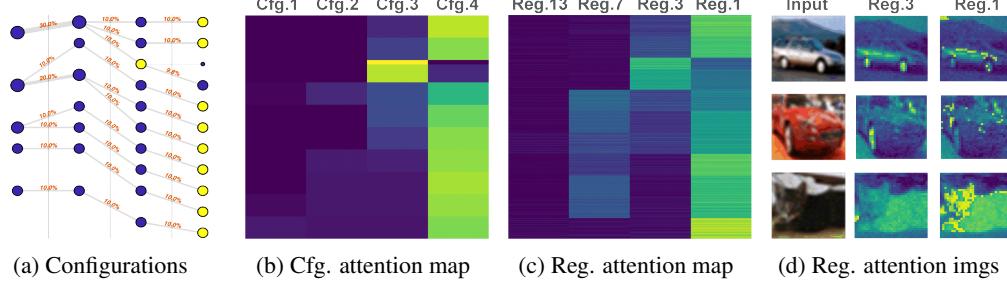


Figure 2: Comparison of attention maps obtained from configurations and registers, rows for samples. **(a)**: Lineage diagram for configurations, near GT balls are marked yellow. **(b)**: Attention map of configuration tokens in an attention-based linear probing. **(c)**: Attention map of DINOv2-reg register tokens, mean of all patch norms is used. **(d)**: Attention maps over the register tokens, as images.

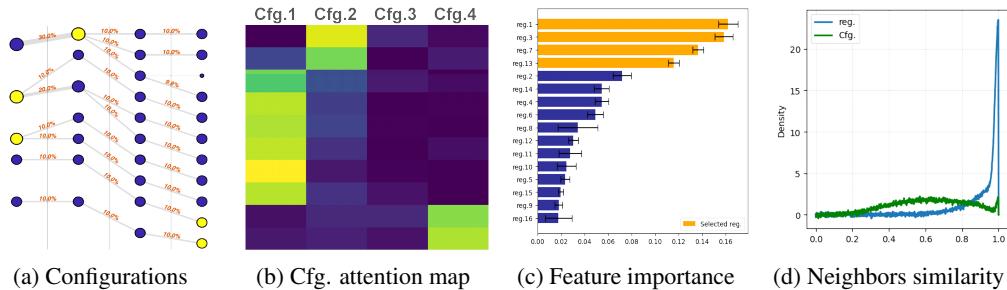


Figure 3: Illustration of another two properties of configurations, grouped by left two and right two. **(a)**: Lineage diagram where coarser classes are used for GT. **(b)**: Attention map in linear probing the coarser classes. **(c)**: Distribution of feature vector importance over the register tokens querying, mean of all patch importance is used. **(d)**: Distribution of cosine similarity between query embeddings of register and configuration tokens and their 2 neighbors, mean of all patch similarities is used.

93 **Configurations are more informative and less redundant than register tokens.** Register tokens  
 94 can help extract configurations, similar to object detection [20, 21], but selecting a fixed number by  
 95 feature importance is arbitrary and non-rigorous (see Fig. 3c). Furthermore, register tokens exhibit  
 96 high redundancy—cosine similarity between their embeddings and their 2 neighbors embeddings is  
 97 heavily skewed toward 1—whereas configurations yield information less redundant (see Fig. 3d).

### 98 3 Methodology

99 Having these characterizations, we hypothesize that unsupervised methods can produce hierarchical  
 100 *multi-resolution clusterings*, and that task- and input-specific *selection and mixing* of these configurations  
 101 represent *global information* beneficial to downstream tasks. Building on the hypothesis, we  
 102 propose a lightweight module *GraMixC*, that treats configurations as tokens ([CFG]) and incorporates  
 103 a novel alignment layer plus learnable attention heads [22] after the configuration extraction model,  
 104 enabling task- and input-specific mixing of configurations via end-to-end back-propagation.

105 Fig. 4 illustrates *GraMixC*. Given an input matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$  (with  $N$  samples and feature dimension  $d$ ), *GraMixC* pass  $\mathbf{X}$  to two branches: (1) a path to unsupervised learning box that extracts  
 106 configurations, and (2) a direct path to the downstream predictor. If at inference, we apply *Reverse*  
 107 *Merge & Split* (RMS) alignment on the configurations. Then we pass them to positional encoding (PE)  
 108 and attention heads. The final concatenation is passed to a downstream predictor for the prediction  $\hat{y}$ .

109 Except for the downstream predictor, the *GraMixC* model can be divided into three parts: the  
 110 unsupervised learning of configurations, the Reverse Merge & Split (RMS) for alignment, and  
 111 attention heads for fusion. In the attention heads part, following Dariset et al. [16], we append register  
 112 tokens [REG] after [CFG] and [CLS] for a clean attention map, that can be used backwards to guide  
 113 configuration selection. Below we detail the rest two components in Section 3.1 and Section 3.2.

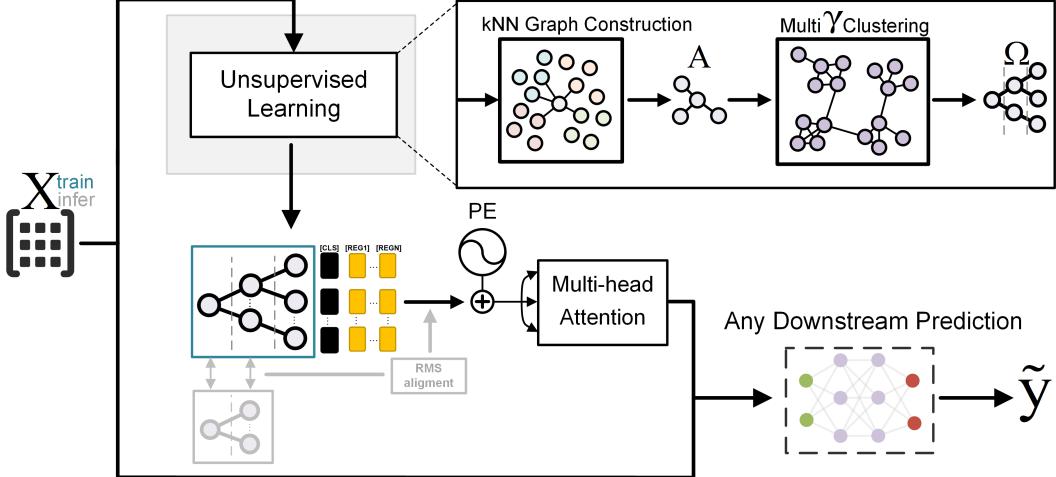


Figure 4: Illustration of the proposed GraMixC module and resulting model. The input data branches into (upper) a path to unsupervised learning box that extracts configurations, and (lower) a direct path to the downstream predictor. Their outcomes concatenate and pass to the downstream predictor. The components occur only during training and inference are colored in blue and gray, respectively.

### 115 3.1 Multi-resolution graph-based clustering

116 Given  $\mathbf{X}$ , multi-resolution clustering seeks to extract *configurations*—valid hierarchical clusterings  
 117 across multiple resolution scales—which we denote as  $\Omega \in \mathbb{N}^{N \times m}$ , where  $m$  denotes the number of  
 118 valid resolution levels. To preserve the latent manifold structure in data, ease parameter sensitivity, and  
 119 prevent other problems with traditional clustering methods (see Section B), we choose the resolution  
 120 parameter ( $\gamma \in \mathbb{R}_+$ )-based community detection as our core clustering method. While BlueRed [23]  
 121 can conduct graph clustering without problems like resolution limit or parameter sensitivity in  
 122 traditional methods, recent work by Pitsianis et al. [8] further demonstrates the elimination of  $\gamma$   
 123 selection, and enabled the unsupervised discovery of  $\Omega$  and the corresponding set of all valid  $\gamma$ , which  
 124 is denoted as  $\Gamma = \{\gamma_1^*, \gamma_2^*, \dots, \gamma_m^*\} \subseteq [0, \infty)$ . Inspired by these works, the unsupervised box in  
 125 Fig. 4 unfolds into two steps: **(1) k-nearest neighbors (kNN) [24] graph construction**, which return  
 126 a directed graph  $G = (V, E)$ , usually represented as adjacency matrix  $\mathbf{A} \in \mathbb{R}_{+}^{N \times N}$ , and **(2) multi- $\gamma$**   
 127 **clustering** on the resulted graph, i.e. modularity based community detection with unsupervised  $\Gamma$   
 128 learning, which return the wanted  $\Omega$ . The details for each of these two steps are:

129 **(1) kNN graph construction.** We construct a kNN graph with  $k = \log_{10} N$  as convention, using  
 130 Euclidean distance for simplicity. Such pair-wise geometric distance between two different vertexes  
 131 is denoted  $d(\mathbf{x}_i, \mathbf{x}_j)$  where  $i \neq j$  and  $\mathbf{x}_i \in \mathbb{R}^d$  is the  $i$ -th feature vector. We then have the adjacency  
 132 matrix  $\mathbf{A}$  formulated as:  $A_{ij} = d(\mathbf{x}_i, \mathbf{x}_j)$  if  $(\mathbf{x}_i, \mathbf{x}_j) \in E$ , 0 otherwise, where  $E$  is the edge set  
 133 of the kNN graph and  $A_{ij}$  denotes the  $i$ -th row and  $j$ -th column element of the adjacency matrix.  
 134 Then we force *column stochastic* by dividing each column in the constructed  $\mathbf{A}$  with the column  
 135 sum. The resulted graph is sparse stochastic, and we can apply Stochastic Graph t-SNE (SG-t-SNE)  
 136 reweighting [25], which proved to remedy skewed degree distribution, that is not promised by  
 137 conventional t-SNE [26]. From the original work, the key equations for SG-t-SNE reweighting are:

$$w(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{\lambda} \exp \left( -\frac{d^2(\mathbf{x}_i, \mathbf{x}_j)}{2\sigma_i^2} \right), \quad \text{with } \lambda = \sum_{\mathbf{x}_j: (\mathbf{x}_i, \mathbf{x}_j) \in E} \exp \left( -\frac{d^2(\mathbf{x}_i, \mathbf{x}_j)}{2\sigma_i^2} \right),$$

138 where  $\lambda$  is a non-negative parameter constant, which we simply set to 15 as previous work show  
 139 that it is not so sensitive to the choice of  $\lambda$  [25], and  $\sigma_i$  is a variable to be numerically solved with  
 140 bisection method. After giving value of  $w$  to  $d$ , we have  $\mathbf{A}$  with less skewed degree distribution,  
 141 which avoids problems like numerical instability and bias towards hubs in downstream clustering.

142 **(2) multi- $\gamma$  community detection.** Then one may simply pass the reweighted  $\mathbf{A}$  to  $\gamma$ -based com-  
 143 munity detection method, such as Leiden algorithm [27], to get one pseudo-configuration vector  
 144  $\omega_\gamma \in \{1, \dots, N\}^N$  (“pseudo” for not sure to be valid). However, such  $\gamma$  falls in the range of  $[0, \infty)$ ,  
 145 and searching over all possible  $\gamma$  is exhausting. Therefore, we incorporate the BlueRed method with

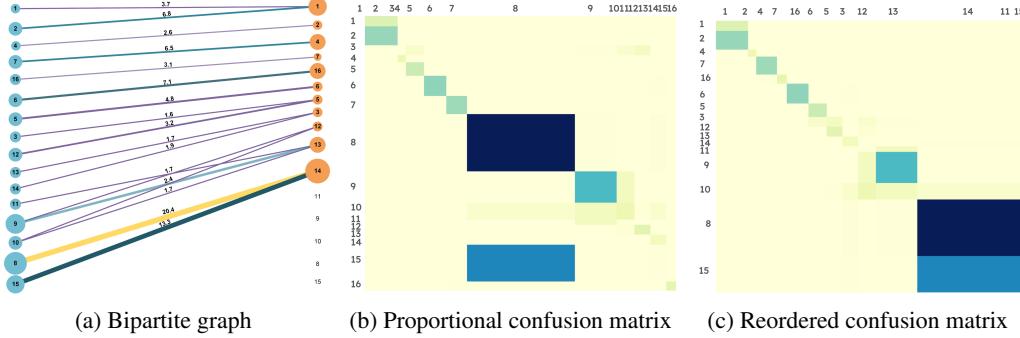


Figure 5: Example of the RMS alignment process applied to clustering results and ground truth (both treated as configurations) on the Salinas dataset [28]. (a): Bipartite graph representation, where blue nodes correspond to predicted clusters and red nodes to ground truth clusters. Node labels indicate cluster indices; edge labels show the proportion of samples shared between clusters. (b): Proportional confusion matrix  $C$  comparing predicted clusters (horizontal axis) to ground truth clusters (vertical axis). (c): Confusion matrix  $C_{tw}$  reordered via the two-walk Laplacian. Notable splits, such as ground truth cluster 8 being divided into clusters 8 and 15 in the prediction, can be resolved through the reverse merge/split procedure.

146 parallel descending triangulation (parallel-DT) [8], in order to automatically discover all valid  $\gamma^* \in \Gamma$ .  
 147 Given a fixed  $\gamma$ , BlueRed find the optimal configuration  $\omega_\gamma$  by the following optimization:

$$\omega_\gamma = \arg \min_{\omega \in \{1, \dots, N\}^N} \left[ - \sum_{k=1}^{|\omega|_\infty} \sum_{(i,j) \in E} d(\mathbf{x}_i, \mathbf{x}_j) \mathbf{1}_{\omega_i = \omega_j = k} + \gamma \sum_{k=1}^{|\omega|_\infty} \sum_{(i,j) \in E} d^2(\mathbf{x}_i, \mathbf{x}_j) \mathbf{1}_{\omega_i = k}, \right],$$

148 where  $\omega_i$  denotes the  $i$ -th element of  $\omega$ ,  $|\omega|_\infty = \max_{i \leq N} \omega_i$  is a inf-norm, and  $\mathbf{1}$  denotes the indicator  
 149 gate which take value 1 if its subscript condition holds, 0 otherwise. Pitsianis et al. [8] describe the  
 150 first term as attraction and the second term as repulsion. Optimizing each solely yields all-in-one  
 151 configuration  $\omega_0 = [1, 1, \dots, 1]$  and all-lonely configuration  $\omega_\infty = [1, 2, \dots, N]$ . Between these  
 152 two configurations, parallel-DT allows forming BlueRed Front (BRF) [8] by segmenting  $(0, \infty)$  into  
 153  $m$  ranges, among which each has a dominant  $\gamma_i^*$  yields lower HAR [8]—the sum of first term and the  
 154 negative second term—which means “local minimum” on that range. Thus desired  $\Omega$  is formed.

### 155 3.2 RMS: reverse merge & split alignment

156 Multi-resolution clustering on different datasets  $\mathbf{X}_{train}$  and  $\mathbf{X}_{test}$  often naturally produces misaligned  
 157 configurations, that either (1) have different value of  $m$  or  $|\omega|_\infty$ , or (2) have different cluster labels.  
 158 While (2) is not a problem as re-assigning fix it, (1) could be problematic as the length and position  
 159 of configurations influence the downstream fusion. One possible interpretation is that some clusters  
 160 are further merged or split in another configuration, leading to this mismatch. To address this, we  
 161 propose Reverse Merge & Split (RMS), which identifies an optimal alignment, allowing re-merging  
 162 and re-splitting, between two configurations,  $\omega_i$  and  $\omega_j$ . First of all, an alignment score is defined:

$$SCORE(\omega_i, \omega_j) = ARI(\omega_i, \omega_j) - \theta \left| \frac{|\omega_i|_\infty - |\omega_j|_\infty}{|\omega_i|_\infty + |\omega_j|_\infty} \right|.$$

163 where  $\theta$  is a hyperparameter to balance the weights of the two terms, which we set to 0.1, ARI is the  
 164 adjusted rand index as defined in Hubert and Arabie [29]. By this punished ARI design, we consider  
 165 different labels, merge and split during scoring the alignment between two partition, but also avoids  
 166 too much difference in number of clusters (one extreme case is  $\omega_0$  and  $\omega_\infty$  has ARI of 1).

167 However, the SCORE itself does not convey the mapping we need for reassigning. In RMS alignment,  
 168 we construct a confusion matrix  $C \in \mathbb{N}^{|\omega_i|_\infty \times |\omega_j|_\infty}$  between  $\omega_i$  and  $\omega_j$ . Fig. 5 illustrates this  
 169 process with a concrete example, showing how the confusion matrix captures the relationship  
 170 between predicted and ground truth clusters, including cases where clusters are split or merged across  
 171 configurations. As an assignment problem with a rectangle cost matrix  $-C$ ,<sup>1</sup> it is solvable by twisting

<sup>1</sup>The negative of the confusion matrix is used to frame the assignment problem (minimizing the diagonal).

existing Hungarian algorithm methods [30, 31, 32]. Because  $C$  is the adjacency matrix of a bipartite graph, spectral reordering via its graph Laplacian is preferred, since it encodes global connectivity and reveals coherent split–merge structures rather than merely optimizing diagonal entries. As the Fiedler vector reordering [33] assumes symmetric positive semi-definite, it is not directly applicable to  $C$ . Inspired by a recent work of Floros et al. [34], we introduce a two-walk Laplacian, which is defined as:

$$L_{tw} = D - C_{tw}, \quad \text{with } C_{tw} = \begin{bmatrix} CC^\top & C \\ C^\top & C^\top C \end{bmatrix},$$

where  $D = \text{diag}(C_{tw}\mathbf{1})$  is the diagonal degree matrix of  $C_{tw}$ . We remap  $\omega_i$  and  $\omega_j$  by using, respectively, the first  $\|\omega_i\|_\infty$  and the last  $\|\omega_j\|_\infty$  entries in the Fiedler eigenvector of  $L_{tw}$ , which is the eigenvector corresponds to smallest positive eigenvalue. We further reverse split and merge simply by reassigning the redundant columns or rows who has element larger than its diagonal entry.

In GraMixC, we carry a small portion (0.1%) of train samples as anchors during inference, and the portion of  $\Omega_{train}$  and  $\Omega_{test}$  corresponding to the anchors are used to calculate the SCORE. Given  $m$  is usually small, we exhaustively test pairs  $(\omega_i, \omega_j)$  then iteratively pick the pair yielding the highest SCORE for each  $\omega_i$ . For each pair, we apply the mapping from RMS( $\omega_i, \omega_j$ ). The final alignments is then used to match the configurations. See the GitHub repository<sup>2</sup> and Section C for alignment examples and more implementation details.

## 4 Experiments

In this section, we evaluate the proposed plug-and-play module by training baseline models with and without GraMixC (GMC). We also test a static variant (GC), which use aligned configurations as extra features, without attention mechanism. We expect the performance to follow a general trend

$$\text{baseline} < \text{baseline+GC} < \text{baseline+GMC}.$$

We then ablate the number of configurations used to check that they cause a performance regression.

### 4.1 Implementation details and experimental setup

Our module was implemented with MATLAB, Python 3.12, PyTorch 2.6. We run trainings on a GeForce RTX 3090Ti GPU. Models were trained with the Adam optimizer [35] at a fixed learning rate of  $10^{-3}$ . Unless otherwise noted, we used a batch size of 100 and trained for up to 100 epochs.

Ahead of diving into the experimental details, we briefly summarize the datasets and metrics used.

**DSNI-pH and DSNI-Temp.** We collected the DSNI dataset from DSMZ [36] and NIH [37]. It comprises six relational tables (STRAINS, MEDIA, SOLUTIONS, INGREDIENTS, STEPS, GAS) covering taxonomic and protocol information. We use approximately 65 000 samples with 16S rRNA sequence (500–1 500 nucleotides), cultivation temperatures (2–103 °C), and pH (0.1–11). The task is to predict optimal temperature (DSNI-Temp) and pH (DSNI-pH) from the 16S rRNA sequence.

Following Çelikkanat et al. [38] and related works [39, 40], we encode each 16S rRNA sequence as a 7-mer count vector in  $\mathbb{N}^{16,384}$ , yielding a dataset of shape  $65\,023 \times 16\,384$ . We perform an 80/20 split (52,018 train / 13,005 test), which preserves the skewed pH (6–8) and temperature (20–40 °C) distributions. Fig. 6 provides an illustration for target value ( $y_{train}$  and  $y_{test}$ ) distribution. Preprocessing—robust scaling, variance thresholding, and selection of the top 1,000 features—was fitted on the training set and then applied to both splits to avoid data leakage.

**Additional benchmarks.** We further evaluate on QM9 [41] for molecular property regression, on Boston Housing [42], and on MNIST [43] and CIFAR10 for classification (some in Section D).

**Evaluation metrics.** For regression we use mean squared error (MSE), mean absolute error (MAE; used for QM9 for comparability with SOTA) for training, and report coefficient of determination ( $R^2$ ). For classification we use cross-entropy loss (CE) for training and report top-1 accuracy (Acc).

For each benchmark, we include three classical decision tree models for reference: Random Forest (RF) [44], XGBoost [45], CatBoost [46]. As both GMC and GC are plug-and-play modules, they can be easily applied to various downstream predictors. We first evaluate a 3-layer perceptron (3LP) with

<sup>2</sup><https://anonymous.4open.science/r/project-82CE>

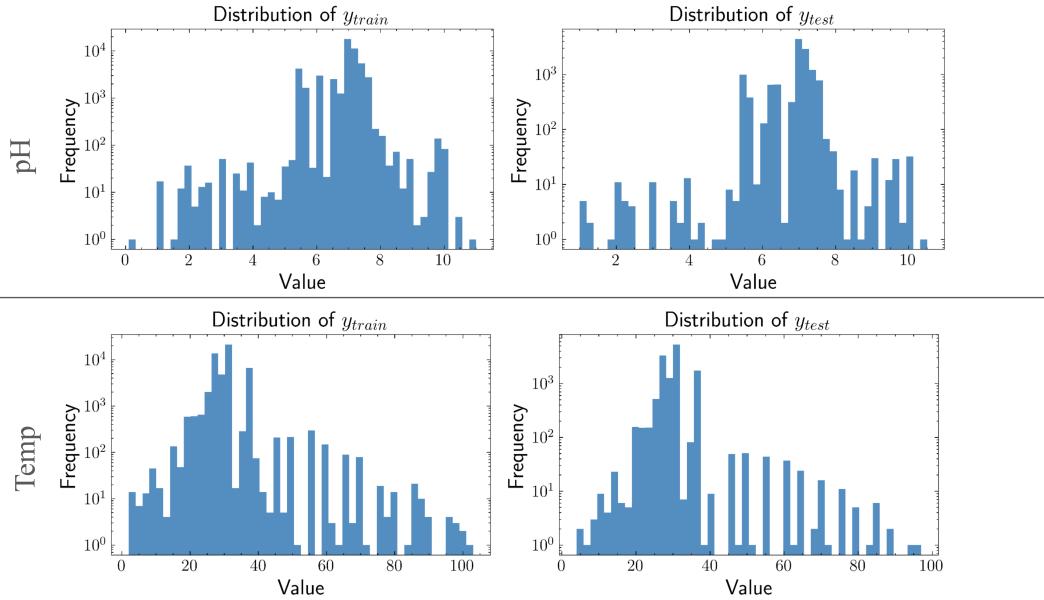


Figure 6: Illustration of target value distributions across train-test splits in DSNI dataset. The first row represents pH distributions and the second row represents temperature distributions. The first column represents the training set ( $y_{train}$ ) and the second column represents the test set ( $y_{test}$ ).

217 hidden dims [256,128,64]. Because our inputs combine numerical features with categorical config-  
 218 urations, we naturally consider tabular models: TabNet (TabN) [47], TabTransformer (TabT) [48],  
 219 FT-Transformer (FTT) [49] were all run with their default settings from the official implementations.

## 220 4.2 Evaluation of the proposed module

221 As shown in Fig. 2 and Fig. 3, we demonstrate, with attention maps, the learned mixing of config-  
 222 urations by training models with self-attention head on aligned configurations. In order to quantify  
 223 the quality of such mixing, for each baseline, we set up the evaluation in three modes: standalone  
 224 (baseline), with static configuration concatenation (baseline+GC), and with attention-based fusion via  
 225 GraMixC (baseline+GMC). Table 1 reports regression results on our main benchmarks; Section D  
 226 (Table 2) shows the rest results. Across all models and tasks, adding GC yields consistent gains, and  
 227 incorporating GMC provides further significant improvements, confirming our initial hypothesis.

228 **Performance improvement.** Table 1 shows that adding GC and GMC yields consistent gains across  
 229 all baselines. Among these observed improvements, the scores increasing on DSNI is quite satisfying.  
 230 Prior specialized growth-media regression methods are not convincing with  $R^2 \leq 0.8$  (e.g., 0.75 [50]).  
 231 We confirm this with our base models score  $R^2$  between 0.3 and 0.6 on DSNI-pH and DSNI-Temp.  
 232 However, even without tailoring the baseline model design, we bring the score to a new high by  
 233 simply adding GC or GMC. Fig. 7 illustrates some examples of such improvement. We see the  
 234 model’s predictions align more closely with the ideal regression line and better handle rare cases,  
 235 by incorporating configurations and probably capturing the latent manifold structure. Incorporating  
 236 GC and further GMC raises  $R^2$  to 0.98 (pH) and 0.97 (Temp). Which not only is considered very  
 237 satisfying in application of bacterial cultivation but also set the new state-of-the-art (SOTA) for  
 238 growth-media prediction. On QM9, GraMixC achieves an MAE of 0.008, nearly matching the SOTA  
 239 (w/o extra training data) of 0.007 [51], and represents the best result among non-GNN models.

240 **Number of configurations used.** We ablate the number of configuration levels in GMC. Fig. 8 shows  
 241 that more configurations generally decreases MSE and increases  $R^2$ , confirming the value of multi-  
 242 resolution information. Importantly, GMC often needs more than half as many total configurations to  
 243 outperform GC, and performance plateaus—or even slightly declines—when including the last few  
 244 configurations. These aligns with Pittsianis et al. [8], who report a finite set of optimal configurations  
 245 rather than continuous gains at infinite resolutions. Using all configurations available is still preferred.

Table 1: Regression performance on DSNI-pH, DSNI-Temp and QM9. Values are mean $\pm$ std from runs with different random seeds; best results per baseline are bold; best results per metric are underlined.

	DSNI-pH		DSNI-Temp		QM9	
	MSE $\downarrow$	R <sup>2</sup>	MSE $\downarrow$	R <sup>2</sup>	MAE $\downarrow$	R <sup>2</sup>
RF	0.198 $\pm$ 0.000	0.601 $\pm$ 0.001	17.759 $\pm$ 0.276	0.393 $\pm$ 0.009	0.015 $\pm$ 0.000	0.979 $\pm$ 0.000
XGBoost	0.196 $\pm$ 0.001	0.604 $\pm$ 0.003	18.212 $\pm$ 0.543	0.377 $\pm$ 0.018	0.014 $\pm$ 0.001	0.978 $\pm$ 0.001
CatBoost	0.193 $\pm$ 0.001	0.610 $\pm$ 0.002	17.375 $\pm$ 0.398	0.406 $\pm$ 0.013	0.014 $\pm$ 0.000	0.978 $\pm$ 0.002
3LP	0.201 $\pm$ 0.002	0.595 $\pm$ 0.006	18.484 $\pm$ 0.183	0.368 $\pm$ 0.006	0.018 $\pm$ 0.001	0.958 $\pm$ 0.001
3LP+GC	0.097 $\pm$ 0.004	0.804 $\pm$ 0.008	6.520 $\pm$ 0.360	0.777 $\pm$ 0.012	0.016 $\pm$ 0.003	0.974 $\pm$ 0.000
3LP+GMC	<b>0.023</b> $\pm$ 0.002	<b>0.953</b> $\pm$ 0.004	<b>2.277</b> $\pm$ 0.061	<b>0.922</b> $\pm$ 0.002	<b>0.010</b> $\pm$ 0.003	<b>0.990</b> $\pm$ 0.002
TabN	0.184 $\pm$ 0.004	0.629 $\pm$ 0.007	13.290 $\pm$ 0.244	0.545 $\pm$ 0.008	0.015 $\pm$ 0.001	0.962 $\pm$ 0.002
TabN+GC	0.086 $\pm$ 0.003	0.825 $\pm$ 0.007	7.997 $\pm$ 0.210	0.726 $\pm$ 0.007	0.012 $\pm$ 0.002	0.983 $\pm$ 0.001
TabN+GMC	<b>0.020</b> $\pm$ 0.001	<b>0.959</b> $\pm$ 0.002	<b>0.989</b> $\pm$ 0.361	<b>0.966</b> $\pm$ 0.012	<b>0.008</b> $\pm$ 0.000	<b>0.995</b> $\pm$ 0.002
TabT	0.256 $\pm$ 0.007	0.483 $\pm$ 0.014	18.910 $\pm$ 0.247	0.353 $\pm$ 0.008	0.434 $\pm$ 0.008	0.921 $\pm$ 0.008
TabT+GC	0.106 $\pm$ 0.002	0.786 $\pm$ 0.005	8.280 $\pm$ 0.303	0.717 $\pm$ 0.010	0.212 $\pm$ 0.004	0.961 $\pm$ 0.008
TabT+GMC	<b>0.017</b> $\pm$ 0.002	<b>0.964</b> $\pm$ 0.005	<b>2.785</b> $\pm$ 0.540	<b>0.904</b> $\pm$ 0.018	<b>0.009</b> $\pm$ 0.000	<b>0.998</b> $\pm$ 0.001
FTT	0.218 $\pm$ 0.003	0.561 $\pm$ 0.006	13.571 $\pm$ 0.069	0.536 $\pm$ 0.002	0.085 $\pm$ 0.005	0.984 $\pm$ 0.006
FTT+GC	0.070 $\pm$ 0.003	0.858 $\pm$ 0.007	5.915 $\pm$ 0.277	0.797 $\pm$ 0.009	0.034 $\pm$ 0.002	0.993 $\pm$ 0.003
FTT+GMC	<b>0.007</b> $\pm$ 0.005	<b>0.984</b> $\pm$ 0.009	<b>1.480</b> $\pm$ 0.120	<b>0.949</b> $\pm$ 0.004	<b>0.026</b> $\pm$ 0.001	<b>0.995</b> $\pm$ 0.003

### 246 4.3 Qualitative evaluation of configurations.

247 Our final experiment compares configurations against standard representation-extraction methods.  
248 As discussed in Section 1, configurations can be viewed as special unsupervised representation  
249 learning. Fig. 3 already shows their advantage over self-supervised register tokens. Here, we replace  
250 GC/GMC with PCA [52], UMAP [53], and a vanilla autoencoder (AE), each embed into dimensions  
251 the same number of as our configurations. We visualize these embeddings on MNIST (Fig. 9a;  
252 additional views in Section D.2). Qualitatively, SG-t-SNE (the reduction step in GraMixC) yields  
253 more uniform, well-separated clusters that respect global kNN connectivity rather than forming hubs.  
254 Fig. 9b quantifies downstream classification accuracy, where GC and GMC strongly outperform PCA,  
255 UMAP, and AE given the same embedding budget. These results confirm that mixed configurations  
256 provide a more expressive yet compact representation for downstream tasks.

## 257 5 Conclusion

258 In this study, we investigate the functional mechanisms of configurations in downstream prediction  
259 tasks and identify three key properties. Based on this, we propose GraMixC, which dynamically  
260 mixes configurations through attention head. We apply it to the challenging task of 16S rRNA  
261 cultivation-media prediction task, and set a new state-of-the-art. Further validation across multiple  
262 standard tabular data benchmarks consistently reveals that GC (a static version of GraMixC) enhances  
263 baseline performance, while GraMixC demonstrates even more substantial improvements. Our results  
264 suggest that harnessing rich manifold priors via attention-driven fusion opens promising avenues for  
265 interpretable and robust learning in both scientific and conventional domains.

266 In future work, we plan to extend mixed configurations to more expressive networks and dynamically  
267 learn configuration alignment through end-to-end differentiable modules. Additionally, we will focus  
268 on exploring adaptive clustering for evolving data streams where train and test distributions may shift,  
269 which could further enhance the resilience of multi-resolution approaches.

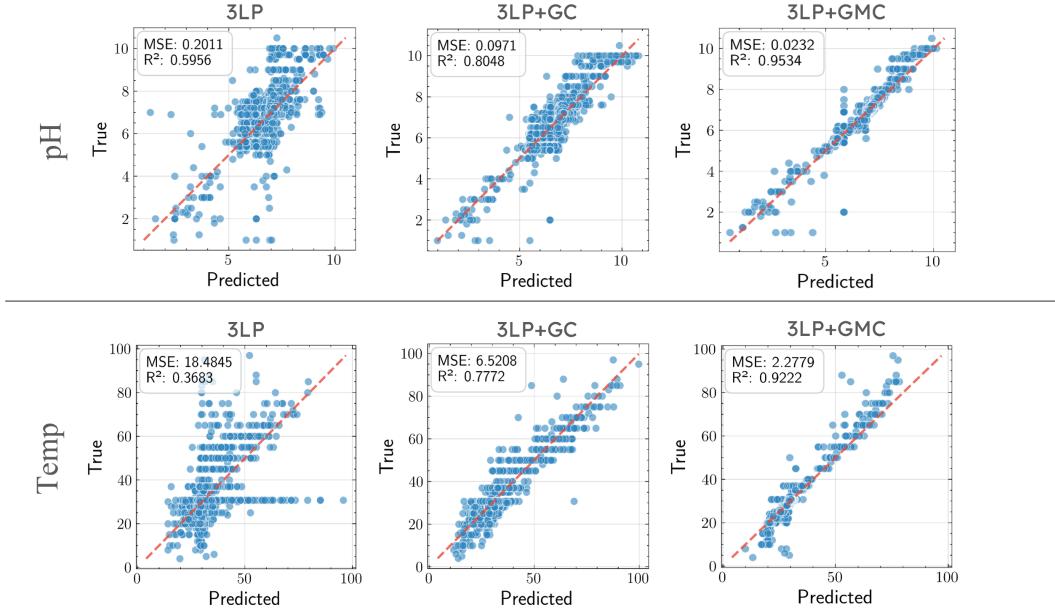


Figure 7: Illustration of the regression performance improvement example in 3LP by adding GC or GMC. Each column plots predicted vs. actual pH (top) or temperature (bottom). 3LP+GC (middle) outperforms the 3LP baseline (left), while 3LP+GMC (right) further boosts  $R^2$  up to  $> 0.9$ .

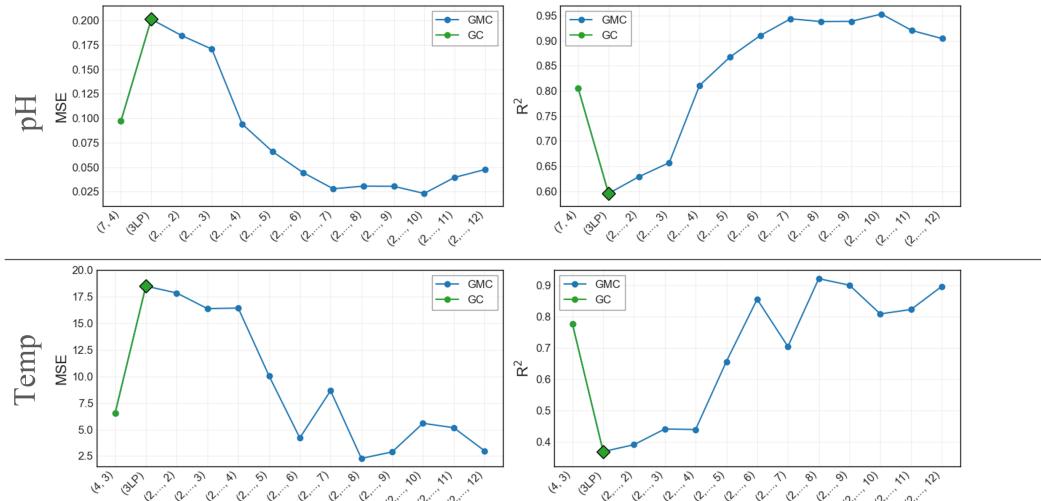
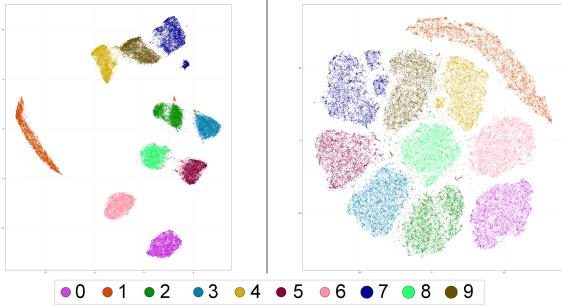


Figure 8: Ablation study on the number of configurations used on DSNI. On the blue curves (GMC),  $[2, \dots, i]$  denote fusing configurations from 2 through  $i$  via GraMixC. On the green curves (GC),  $(i, j)$  denote the best train/test configuration pair used in static concatenation. Incrementally mixing configurations improves performance and outperforms static concatenation.

## 270 References

- 271 [1] J. MacQueen. “Some Methods for Classification and Analysis of Multivariate Observations”.  
272 In: *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*. Vol. 5.1. University of California Press, Jan. 1, 1967, pp. 281–298.
- 274 [2] Jianbo Shi and J. Malik. “Normalized Cuts and Image Segmentation”. In: *IEEE Trans. Pattern  
275 Anal. Machine Intell.* 22.8 (Aug. 2000), pp. 888–905.
- 276 [3] A. Ng, M. Jordan, and Y. Weiss. “On Spectral Clustering: Analysis and an Algorithm”. In:  
277 *Advances in Neural Information Processing Systems*. Vol. 14. MIT Press, 2001.



(a) 2D visualization of embeddings learned.

	CE ↓	Acc
3LP+PCA	0.157	0.971
3LP+UMAP	0.181	0.975
3LP+AE	0.158	0.969
3LP+GC	0.046	0.992
<b>3LP+GMC</b>	<b>0.028</b>	<b>0.993</b>

(b) Classification performance.

Figure 9: (a): Illustration of 2D embeddings of MNIST using UMAP (left) and SG-t-SNE (right). (b): Classification performance on MNIST using features from PCA, UMAP, autoencoder (AE), static configurations (GC), and GraMixC (GMC) at equal embedding dimensions. SG-t-SNE embeddings integrated via GC or GMC exploit multi-resolution structure to notably outperform other methods.

- [4] P. C. Quinn and P. D. Eimas. “Perceptual Cues That Permit Categorical Differentiation of Animal Species by Infants”. In: *J Exp Child Psychol* 63.1 (Oct. 1996), pp. 189–211. PMID: 8812045.
- [5] M. H. Bornstein, M. E. Arterberry, and C. Mash. “Infant Object Categorization Transcends Diverse Object-Context Relations”. In: *Infant Behav Dev* 33.1 (Feb. 2010), pp. 7–15. PMID: 20031232.
- [6] L. Zaadnoordijk, T. R. Besold, and R. Cusack. “Lessons from Infant Learning for Unsupervised Machine Learning”. In: *Nat Mach Intell* 4.6 (June 2022), pp. 510–520.
- [7] L. Muttenhaler, K. Greff, F. Born, B. Spitzer, S. Kornblith, M. C. Mozer, K.-R. Müller, T. Unterthiner, and A. K. Lampinen. *Aligning Machine and Human Visual Representations across Abstraction Levels*. Oct. 29, 2024. arXiv: 2409 .06509 [cs]. URL: <http://arxiv.org/abs/2409.06509> (visited on 05/11/2025). Pre-published.
- [8] N. Pitsianis, D. Floros, T. Liu, and X. Sun. “Parallel Clustering with Resolution Variation”. In: *2023 IEEE High Performance Extreme Computing Conference (HPEC)*. 2023 IEEE High Performance Extreme Computing Conference (HPEC). Boston, MA, USA: IEEE, Sept. 25, 2023, pp. 1–8.
- [9] A. Krizhevsky. *Learning Multiple Layers of Features from Tiny Images*. Toronto, ON, Canada, 2009, pp. 32–33.
- [10] J. M. Janda and S. L. Abbott. “16S rRNA Gene Sequencing for Bacterial Identification in the Diagnostic Laboratory: Pluses, Perils, and Pitfalls”. In: *J Clin Microbiol* 45.9 (Sept. 2007), pp. 2761–2764. PMID: 17626177.
- [11] Q. Wang, G. M. Garrity, J. M. Tiedje, and J. R. Cole. “Naive Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy”. In: *Appl Environ Microbiol* 73.16 (Aug. 2007), pp. 5261–5267. PMID: 17586664.
- [12] J. De Vrieze, A. J. Pinto, W. T. Sloan, and U. Z. Ijaz. “The Active Microbial Community More Accurately Reflects the Anaerobic Digestion Process: 16S rRNA (Gene) Sequencing as a Predictive Tool”. In: *Microbiome* 6.1 (Apr. 2, 2018), p. 63.
- [13] J. S. Johnson, D. J. Spakowicz, B.-Y. Hong, L. M. Petersen, P. Demkowicz, L. Chen, S. R. Leopold, B. M. Hanson, H. O. Agresta, M. Gerstein, E. Sodergren, and G. M. Weinstock. “Evaluation of 16S rRNA Gene Sequencing for Species and Strain-Level Microbiome Analysis”. In: *Nat Commun* 10.1 (Nov. 6, 2019), p. 5029.
- [14] M. Caron, P. Bojanowski, A. Joulin, and M. Douze. *Deep Clustering for Unsupervised Learning of Visual Features*. Version 2. Mar. 18, 2019. arXiv: 1807 .05520 [cs]. URL: <http://arxiv.org/abs/1807.05520> (visited on 05/13/2025). Pre-published.
- [15] Y. Yang, Z. Guan, Z. Wang, W. Zhao, C. Xu, W. Lu, and J. Huang. *Self-Supervised Heterogeneous Graph Pre-training Based on Structural Clustering*. Apr. 12, 2023. arXiv: 2210 .10462 [cs]. URL: <http://arxiv.org/abs/2210.10462> (visited on 05/13/2025). Pre-published.

- 315 [16] T. Dariset, M. Oquab, J. Mairal, and P. Bojanowski. *Vision Transformers Need Registers*.  
 316 Apr. 12, 2024. arXiv: 2309 . 16588 [cs]. URL: <http://arxiv.org/abs/2309.16588>  
 317 (visited on 03/25/2025). Pre-published.
- 318 [17] F. Locatello, D. Weissenborn, T. Unterthiner, A. Mahendran, G. Heigold, J. Uszkoreit, A.  
 319 Dosovitskiy, and T. Kipf. *Object-Centric Learning with Slot Attention*. Oct. 14, 2020. arXiv:  
 320 2006 . 15055 [cs]. URL: <http://arxiv.org/abs/2006.15055> (visited on 05/11/2025).  
 321 Pre-published.
- 322 [18] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. *Emerging  
 323 Properties in Self-Supervised Vision Transformers*. May 24, 2021. arXiv: 2104 . 14294 [cs].  
 324 URL: <http://arxiv.org/abs/2104.14294> (visited on 03/25/2025). Pre-published.
- 325 [19] M. Oquab, T. Dariset, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D.  
 326 Haziza, F. Massa, A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. Huang,  
 327 S.-W. Li, I. Misra, M. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jegou, J. Mairal, P. Labatut,  
 328 A. Joulin, and P. Bojanowski. *DINOv2: Learning Robust Visual Features without Supervision*.  
 329 Feb. 2, 2024. arXiv: 2304 . 07193 [cs]. URL: <http://arxiv.org/abs/2304.07193>  
 330 (visited on 03/25/2025). Pre-published.
- 331 [20] O. Siméoni, G. Puy, H. V. Vo, S. Roburin, S. Gidaris, A. Bursuc, P. Pérez, R. Marlet, and  
 332 J. Ponce. *Localizing Objects with Self-Supervised Transformers and No Labels*. Sept. 29,  
 333 2021. arXiv: 2109 . 14279 [cs]. URL: <http://arxiv.org/abs/2109.14279> (visited on  
 334 03/25/2025). Pre-published.
- 335 [21] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. M. Ni, and H.-Y. Shum. *DINO: DETR with  
 336 Improved DeNoising Anchor Boxes for End-to-End Object Detection*. July 11, 2022. arXiv:  
 337 2203 . 03605 [cs]. URL: <http://arxiv.org/abs/2203.03605> (visited on 03/25/2025).  
 338 Pre-published.
- 339 [22] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Łukasz Kaiser, and  
 340 I. Polosukhin. “Attention Is All You Need”. In: *Advances in Neural Information Processing  
 341 Systems*. Vol. 30. Curran Associates, Inc., 2017.
- 342 [23] T. Liu, D. Floros, N. Pitsianis, and X. Sun. “Digraph Clustering by the BlueRed Method”. In:  
 343 *2021 IEEE High Performance Extreme Computing Conference (HPEC)*. 2021 IEEE High  
 344 Performance Extreme Computing Conference (HPEC). Sept. 2021, pp. 1–7.
- 345 [24] J. B. Tenenbaum, V. de Silva, and J. C. Langford. “A Global Geometric Framework for  
 346 Nonlinear Dimensionality Reduction”. In: *Science* 290.5500 (Dec. 22, 2000), pp. 2319–2323.
- 347 [25] N. Pitsianis, A.-S. Iliopoulos, D. Floros, and X. Sun. “Spaceland Embedding of Sparse  
 348 Stochastic Graphs”. In: *2019 IEEE High Performance Extreme Computing Conference (HPEC)*.  
 349 2019 IEEE High Performance Extreme Computing Conference (HPEC). Sept. 2019, pp. 1–8.
- 350 [26] L. Van der Maaten and G. Hinton. “Visualizing Data Using T-SNE.” In: *Journal of machine  
 351 learning research* 9.11 (2008), pp. 2579–2605.
- 352 [27] V. A. Traag, L. Waltman, and N. J. Van Eck. “From Louvain to Leiden: Guaranteeing Well-  
 353 Connected Communities”. In: *Sci. Rep.* 9.1 (Mar. 26, 2019), p. 5233.
- 354 [28] A. Plaza and J. Tilton. “Automated Selection of Results in Hierarchical Segmentations of Re-  
 355 motely Sensed Hyperspectral Images”. In: *Proceedings. 2005 IEEE International Geoscience  
 356 and Remote Sensing Symposium, 2005. IGARSS ’05*. . 2005 IEEE International Geoscience  
 357 and Remote Sensing Symposium, 2005. IGARSS ’05. Vol. 7. July 2005, pp. 4946–4949.
- 358 [29] L. J. Hubert and P. Arabie. “Comparing Partitions”. In: *Journal of Classification* 2.2–3 (1985),  
 359 pp. 193–218.
- 360 [30] H. W. Kuhn. “The Hungarian Method for the Assignment Problem”. In: *Nav. Res. Logist. Q.*  
 361 2.1–2 (Mar. 1955), pp. 83–97.
- 362 [31] R. Jonker and A. Volgenant. “A Shortest Augmenting Path Algorithm for Dense and Sparse  
 363 Linear Assignment Problems”. In: *Computing* 38.4 (Dec. 1, 1987), pp. 325–340.
- 364 [32] D. P. Bertsekas. “Auction Algorithms for Network Flow Problems: A Tutorial Introduction”.  
 365 In: *Comput Optim Applic* 1.1 (Oct. 1992), pp. 7–66.
- 366 [33] M. Fiedler. “Algebraic Connectivity of Graphs”. In: *Czech. Math. J.* 23.2 (1973), pp. 298–305.
- 367 [34] D. Floros, N. Pitsianis, and X. Sun. “Algebraic Vertex Ordering of a Sparse Graph for Adja-  
 368 cency Access Locality and Graph Compression”. In: *2024 IEEE High Performance Extreme  
 369 Computing Conference (HPEC)*. 2024 IEEE High Performance Extreme Computing Confer-  
 370 ence (HPEC). Sept. 2024, pp. 1–7.

- 371 [35] D. P. Kingma and J. Ba. “Adam: A Method for Stochastic Optimization”. In: *arXiv:1412.6980*  
 372 [*cs.LG*] (Jan. 30, 2017). arXiv: 1412.6980 [*cs.LG*].
- 373 [36] German Collection of Microorganisms and Cell Cultures GmbH. *DSMZ – German Collection*  
 374 *of Microorganisms and Cell Cultures*. 2025.
- 375 [37] *National Institutes of Health (NIH)*. National Institutes of Health (NIH). URL: <https://www.nih.gov/> (visited on 02/04/2025).
- 377 [38] A. Çelikkannat, A. R. Masegosa, and T. D. Nielsen. “Revisiting K-mer Profile for Effective and  
 378 Scalable Genome Representation Learning”. In: () .
- 379 [39] *Kraken: Ultrafast Metagenomic Sequence Classification Using Exact Alignments | Genome*  
 380 *Biology | Full Text*. URL: <https://genomebiology.biomedcentral.com/articles/10.1186/gb-2014-15-3-r46> (visited on 02/04/2025).
- 382 [40] *How to Apply de Bruijn Graphs to Genome Assembly | Nature Biotechnology*. URL: <https://www.nature.com/articles/nbt.2023> (visited on 02/04/2025).
- 384 [41] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. von Lilienfeld. “Quantum Chemistry  
 385 Structures and Properties of 134 Kilo Molecules”. In: *Sci Data* 1.1 (Aug. 5, 2014), p. 140022.
- 386 [42] D. Harrison and D. L. Rubinfeld. “Hedonic Housing Prices and the Demand for Clean Air”.  
 387 In: *Journal of Environmental Economics and Management* 5.1 (Mar. 1, 1978), pp. 81–102.
- 388 [43] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. “Gradient-Based Learning Applied to Docu-  
 389 ment Recognition”. In: *Proceedings of the IEEE* 86.11 (Nov. 1998), pp. 2278–2324.
- 390 [44] L. Breiman. “Random Forests”. In: *Machine Learning* 45.1 (2001), pp. 5–32.
- 391 [45] T. Chen and C. Guestrin. “XGBoost: A Scalable Tree Boosting System”. In: *Proceedings of*  
 392 *the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.  
 393 Aug. 13, 2016, pp. 785–794. arXiv: 1603.02754 [*cs.LG*].
- 394 [46] L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Dorogush, and A. Gulin. “CatBoost: Unbiased  
 395 Boosting with Categorical Features”. In: *Advances in Neural Information Processing Systems*.  
 396 Vol. 31. Curran Associates, Inc., 2018.
- 397 [47] S. O. Arik and T. Pfister. *TabNet: Attentive Interpretable Tabular Learning*. Dec. 9, 2020. arXiv:  
 398 1908.07442 [*cs*]. URL: <http://arxiv.org/abs/1908.07442> (visited on 05/14/2025).  
 399 Pre-published.
- 400 [48] X. Huang, A. Khetan, M. Cvitkovic, and Z. Karnin. *TabTransformer: Tabular Data Modeling*  
 401 *Using Contextual Embeddings*. Dec. 11, 2020. arXiv: 2012.06678 [*cs*]. URL: <http://arxiv.org/abs/2012.06678> (visited on 05/05/2025). Pre-published.
- 403 [49] Y. Gorishniy, I. Rubachev, V. Khrulkov, and A. Babenko. *Revisiting Deep Learning Models*  
 404 *for Tabular Data*. Oct. 26, 2023. arXiv: 2106.11959 [*cs*]. URL: <http://arxiv.org/abs/2106.11959> (visited on 05/14/2025). Pre-published.
- 406 [50] D. B. Sauer and D.-N. Wang. “Predicting the Optimal Growth Temperatures of Prokaryotes  
 407 Using Only Genome Derived Features”. In: *Bioinformatics* 35.18 (Sept. 15, 2019), pp. 3224–  
 408 3231.
- 409 [51] X. Fang, L. Liu, J. Lei, D. He, S. Zhang, J. Zhou, F. Wang, H. Wu, and H. Wang. “Geometry-  
 410 Enhanced Molecular Representation Learning for Property Prediction”. In: *Nat Mach Intell*  
 411 4.2 (Feb. 2022), pp. 127–134.
- 412 [52] *Principal Component Analysis*. Springer Series in Statistics. New York: Springer-Verlag, 2002.
- 413 [53] L. McInnes, J. Healy, and J. Melville. “UMAP: Uniform Manifold Approximation and Pro-  
 414 jection for Dimension Reduction”. In: *arXiv:1802.03426 [stat.ML]* (Sept. 18, 2020). arXiv:  
 415 1802.03426 [*stat.ML*].

## 416 A An Intuitive Example of Configuration Mixing

417 To illustrate the necessity of fusing valid clusterings across resolution scales, we use two synthetic  
 418 point-cloud datasets from scikit-learn: “Moons” and “Blobs.” The Blobs dataset is tuned so that no  
 419 single clustering resolution recovers all three clusters. Fig. 10 visualizes each dataset in 3D, using the  
 420 third axis to encode cluster assignments for the corresponding configuration: coarser configuration (1)  
 421 and finer configuration (2). Configuration (1), by lifting some dots above the plane, cleanly separates  
 422 the two Moon arcs but merges two (purple and green) of the Blobs clusters. Configuration (2), by  
 423 itself, fails the Blobs with a different merge (blue and green). Only by fusing both configurations

424 can all clusters be disentangled—the purple dots in (1) that falls down in (2), emerges correct as the  
 425 green cluster. This toy example shows that multi-resolution clusterings alone are insufficient without  
 426 a fusion mechanism. Our GraMixC use attention-based fusion to integrate these scales. While just  
 427 one demonstration, it highlights the broader advantage of mixing configurations in complex settings.

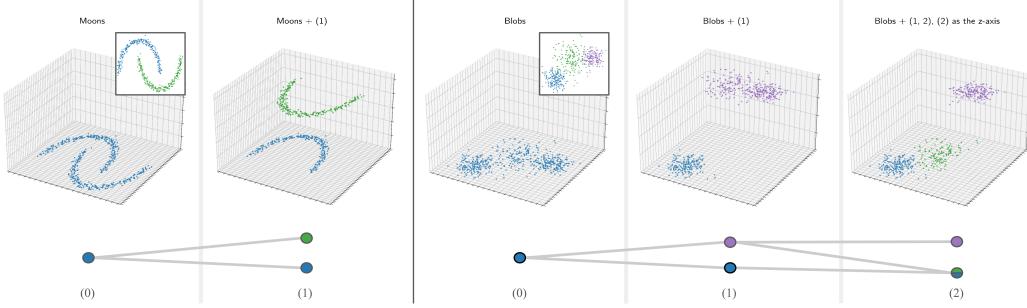


Figure 10: Illustration of multi-resolution clustering on synthetic datasets. GT is shown in the framed box in (0). Upper is the embedding of Moons (left) and Blobs (right) with corresponding configuration ( $i$ ) as third dimension; lower is the lineage diagram of the configurations.

## 428 B Synthetic Clustering Benchmarks

429 In this section, we further discuss the limitations of conventional clustering methods raised in  
 430 Section 3.1. We compare our modularity-based clustering strategy, which is used as the unsupervised  
 431 layer in GraMixC, against widely-used clustering algorithms on synthetic 2D datasets.

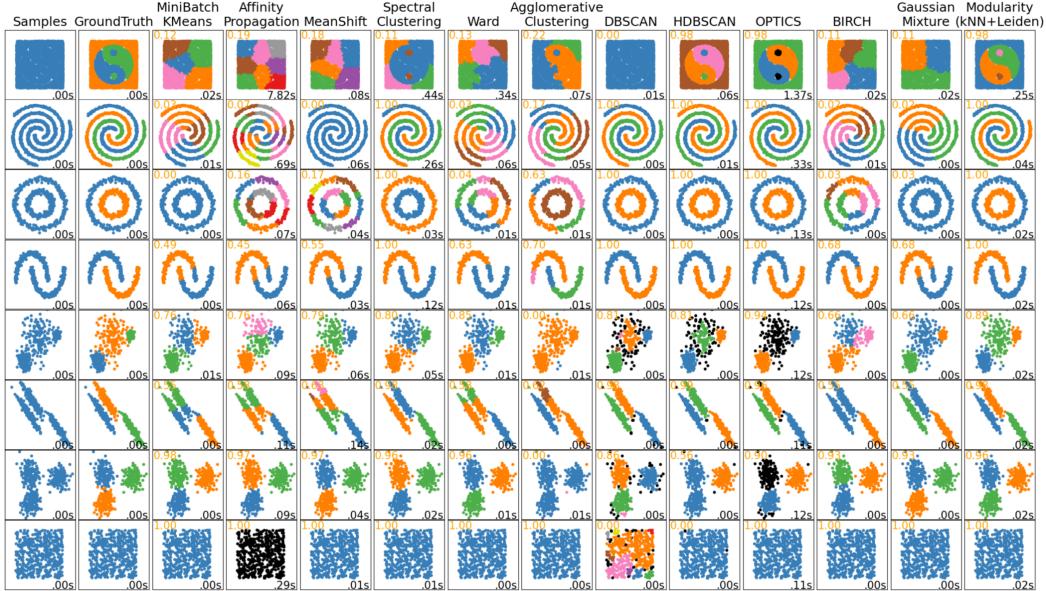


Figure 11: Illustration of clustering methods comparisons across multiple synthetic datasets. Rows correspond to different 2D point clouds—the first row is custom, others are from scikit-learn. Each method’s result is labeled with ARI (top-left in yellow) and execution time (bottom-right in black). **Modularity:** *kNN+Leiden* (far right) accurately recovers ground-truth structures across different shapes and densities, with robustness to noise, anisotropy, and distribution variation.

432 Each row in Fig. 11 presents a distinct synthetic dataset distribution, ranging from custom-designed  
 433 to standard scikit-learn datasets, including *Taiji*, spirals, circles, moons, varied blobs, anisotropic  
 434 blobs, and isotropic noise. Each column represents the result of one clustering method, annotated  
 435 with Adjusted Rand Index (ARI) and execution time.

436 Unlike traditional clustering methods, the approach we adopted (last column: Modularity, im-  
 437 plemented via kNN graph + Leiden community detection) consistently uncovers the underlying  
 438 structure—even in challenging cases involving non-convex geometries, anisotropic spreads, or un-  
 439 even density distributions. This comparison underscores the reliability and manifold sensitivity of our  
 440 unsupervised segmentation approach, even before introducing multi-resolution fusion or downstream  
 441 learning tasks.

## 442 C RMS Alignment Details

443 In Section 3.2 we introduced the Reverse Merge & Split (RMS) procedure for aligning multi-  
 444 resolution configurations between train and test sets. Below we provide the full pseudo-code in  
 445 Algorithm 1, using the same notation as the main text.

### 446 Implementation notes.

- 447 • We set  $\theta = 0.1$  and compute ARI as in Hubert and Arabie [29].
- 448 • We use 0.1 % of the train samples as anchors to form  $\mathcal{A}$ .
- 449 • The greedy matching loops over each train configuration  $\omega_i$  to find its best-scoring test  
 450 partner  $\omega_j$ , applies the label mapping, and removes both from further consideration to ensure  
 451 one-to-one alignment.

452 The details for SCORE and  $L_{tw}$  are covered in Algorithm 1 so we skip them here.

---

#### Algorithm 1 Reverse Merge & Split (RMS) Alignment

---

**Require:**  $\Omega_{\text{train}} \in \mathbb{N}^{N \times m_t}$ ,  $\Omega_{\text{test}} \in \mathbb{N}^{N \times m_s}$ , anchor indices  $\mathcal{A} \subset \{1, \dots, N\}$ ,  $\theta$   
**Ensure:** Aligned  $\Omega_{\text{test}}$

```

1:  $\mathbb{U} \leftarrow \{1, \dots, m_t\}$ ,  $\mathbb{V} \leftarrow \{1, \dots, m_s\}$ 
2: for  $i$  in  $\mathbb{U}$  do ▷ for each train configuration  $\omega_i$ 
3:    $\text{best\_score} \leftarrow -\infty$ ,  $\text{best\_j} \leftarrow \text{null}$ 
4:    $\omega_i \leftarrow \Omega_{\text{train}}[\mathcal{A}, i]$  ▷ find best test configuration  $\omega_j$ 
5:   for  $j$  in  $\mathbb{V}$  do
6:      $\omega_j \leftarrow \Omega_{\text{test}}[\mathcal{A}, j]$ 
7:      $s \leftarrow \text{SCORE}(\omega_i, \omega_j, \theta)$ 
8:     if  $s > \text{best\_score}$  then
9:        $\text{best\_score} \leftarrow s$ ,  $\text{best\_j} \leftarrow j$ 
10:    end if
11:   end for
12:    $M \leftarrow \text{PAIR\_MAPPING}(\Omega_{\text{train}}[:, i], \Omega_{\text{test}}[:, \text{best\_j}])$ 
13:   for  $p = 1$  to  $N$  do
14:      $\Omega_{\text{test}}[p, \text{best\_j}] \leftarrow M(\Omega_{\text{test}}[p, \text{best\_j}])$ 
15:   end for
16:   Remove  $i$  from  $\mathbb{U}$ , remove  $\text{best\_j}$  from  $\mathbb{V}$ 
17: end for
18: return  $\Omega_{\text{test}}$ 

19: function PAIR_MAPPING( $\omega_i, \omega_j$ )
20:    $n_i \leftarrow \|\omega_i\|_\infty$ ,  $n_j \leftarrow \|\omega_j\|_\infty$  ▷ build confusion matrix  $C \in \mathbb{N}^{n_i \times n_j}$ 
21:   for  $p = 1$  to  $N$  do
22:      $C[\omega_i[p], \omega_j[p]] += 1$ 
23:   end for
24:   Construct two-walk Laplacian  $L_{tw}$ 
25:    $\mathcal{F} \leftarrow \text{Fiedler vector of } L_{tw}$ 
26:   Split  $\mathcal{F} \rightarrow (\mathcal{F}_i \in \mathbb{R}^{n_i}, \mathcal{F}_j \in \mathbb{R}^{n_j})$ 
27:    $\pi_i \leftarrow \text{argsort}(\mathcal{F}_i)$ ,  $\pi_j \leftarrow \text{argsort}(\mathcal{F}_j)$ 
28:   return mapping  $k \mapsto \pi_i[\pi_j^{-1}(k)]$  for  $k = 1, \dots, \min(n_i, n_j)$ 
29: end function

```

---

Table 2: Regression/classification performance on Boston Housing (BHouse), MNIST, and CIFAR10.

Dataset	BHouse		MNIST		CIFAR10	
	Metric	MSE ↓	R <sup>2</sup>	CE ↓	Acc	CE ↓
RF	0.022	0.884	0.247	0.969	1.681	0.463
XGBoost	0.022	0.881	0.066	0.980	1.296	0.539
CatBoost	0.016	0.913	0.096	0.975	1.230	0.567
3LP	0.023	0.879	0.141	0.970	1.428	0.524
3LP+GC	0.022	0.882	0.046	0.992	0.480	0.844
3LP+GMC	<b>0.017</b>	<b>0.909</b>	<b>0.028</b>	<b>0.993</b>	<b>0.220</b>	<b>0.949</b>
TabN	0.033	0.822	0.130	0.964	1.499	0.463
TabN+GC	0.021	0.888	0.225	0.941	0.377	0.876
TabN+GMC	<b>0.012</b>	<b>0.936</b>	<b>0.017</b>	<b>0.995</b>	<b>0.077</b>	<b>0.978</b>
TabT	0.035	0.811	0.192	0.980	1.028	0.706
TabT+GC	<b>0.035</b>	<b>0.813</b>	0.040	0.993	1.049	0.704
TabT+GMC	0.061	0.671	<b>0.018</b>	<b>0.994</b>	<b>0.458</b>	<b>0.911</b>
FTT	0.032	0.826	0.098	0.980	0.415	0.874
FTT+GC	0.030	0.838	0.029	0.993	0.437	0.870
FTT+GMC	<b>0.026</b>	<b>0.860</b>	<b>0.018</b>	<b>0.995</b>	<b>0.157</b>	<b>0.955</b>

## 453 D Additional Experimental Results

454 In Section 4 we introduced our experimental setup and high-level results. Here, we provide the full  
455 details and qualitative analyses that couldn’t fit into the main body, including:

- 456 • Downstream task performance on three other benchmarks.  
457 • Qualitative illustration of prediction versus true value on the three tabular baseline models.  
458 • Embeddings from PCA and AE.

### 459 D.1 Additional evaluation of proposed module

460 Table 2 extends our evaluation to three additional benchmarks: Boston Housing (regression), MNIST  
461 and CIFAR-10 (classification). We compare classical ensembles (RF, XGBoost, CatBoost), a 3-layer  
462 MLP (3LP), and three neural tabular architectures (TabNet, TabTransformer, FT-Transformer) in  
463 three modes: baseline, static configuration concatenation (GC), and attention-based fusion (GMC).

464 Across almost all models and datasets, GC consistently improves performance over the raw baselines,  
465 and GMC provides further gains.

466 The sole exception is TabTransformer on Boston Housing, where GC yields only a marginal R<sup>2</sup>  
467 increase (0.811→0.813), but GMC degrades it (to 0.671), suggesting that attention-based fusion may  
468 disrupt already well-structured features in this case.

469 On MNIST, GC lifts accuracy above 99%, and GMC pushes it to 99.3–99.5%. On CIFAR-10, GC  
470 delivers dramatic gains (e.g. TabTransformer from 46.3% to 87.6%), and GMC further improves  
471 all models, with FT-Transformer+GMC reaching 95.5% accuracy. These results underscore that  
472 configuration integration via GraMixC is broadly effective, with only one minor counterexample.

### 473 D.2 Additional qualitative evaluation of configurations

474 In Section 4.3 we provided the embedding of MNIST digits using UMAP and SG-t-SNE (Fig. 9a).  
475 Here we provide the missing illustration of embedding with PCA and autoencoder (AE) in Fig. 12.  
476 As expected, they do not provide representation with clusters as separated as the former two methods.

477 With the final figure (Fig. 13) we visualize predicted vs. actual values from the tabular baselines on  
478 DSNI, filling in what is missing from Fig. 7.

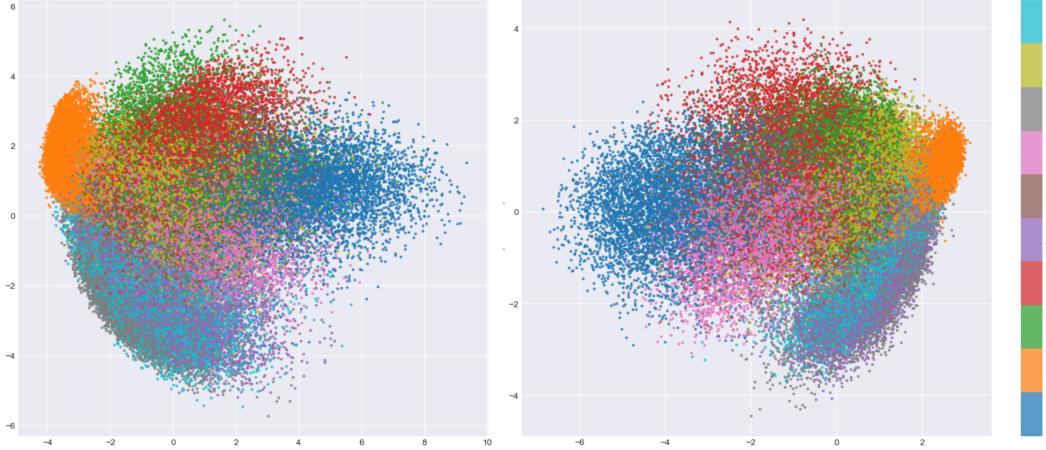


Figure 12: Illustration of 2D embeddings learned by PCA (left) and AE (right) on MNIST.

## 479 NeurIPS Paper Checklist

### 480 1. Claims

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

483 Answer: [Yes]

484 Justification: In the abstract and introduction, we outline configuration characteristics and  
 485 propose GraMixC. We detail our observation of configurations in Section 2 and methods in  
 486 Section 3.

487 Guidelines:

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

### 497 2. Limitations

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

499 Answer: [Yes]

500 Justification: We discuss present limitations and future plans in Section 5.

501 Guidelines:

- 502 • The answer NA means that the paper has no limitation while the answer No means that  
 503 the paper has limitations, but those are not discussed in the paper.
- 504 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 505 • The paper should point out any strong assumptions and how robust the results are to  
 506 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
 507 model well-specification, asymptotic approximations only holding locally). The authors  
 508 should reflect on how these assumptions might be violated in practice and what the  
 509 implications would be.
- 510 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
 511 only tested on a few datasets or with a few runs. In general, empirical results often  
 512 depend on implicit assumptions, which should be articulated.

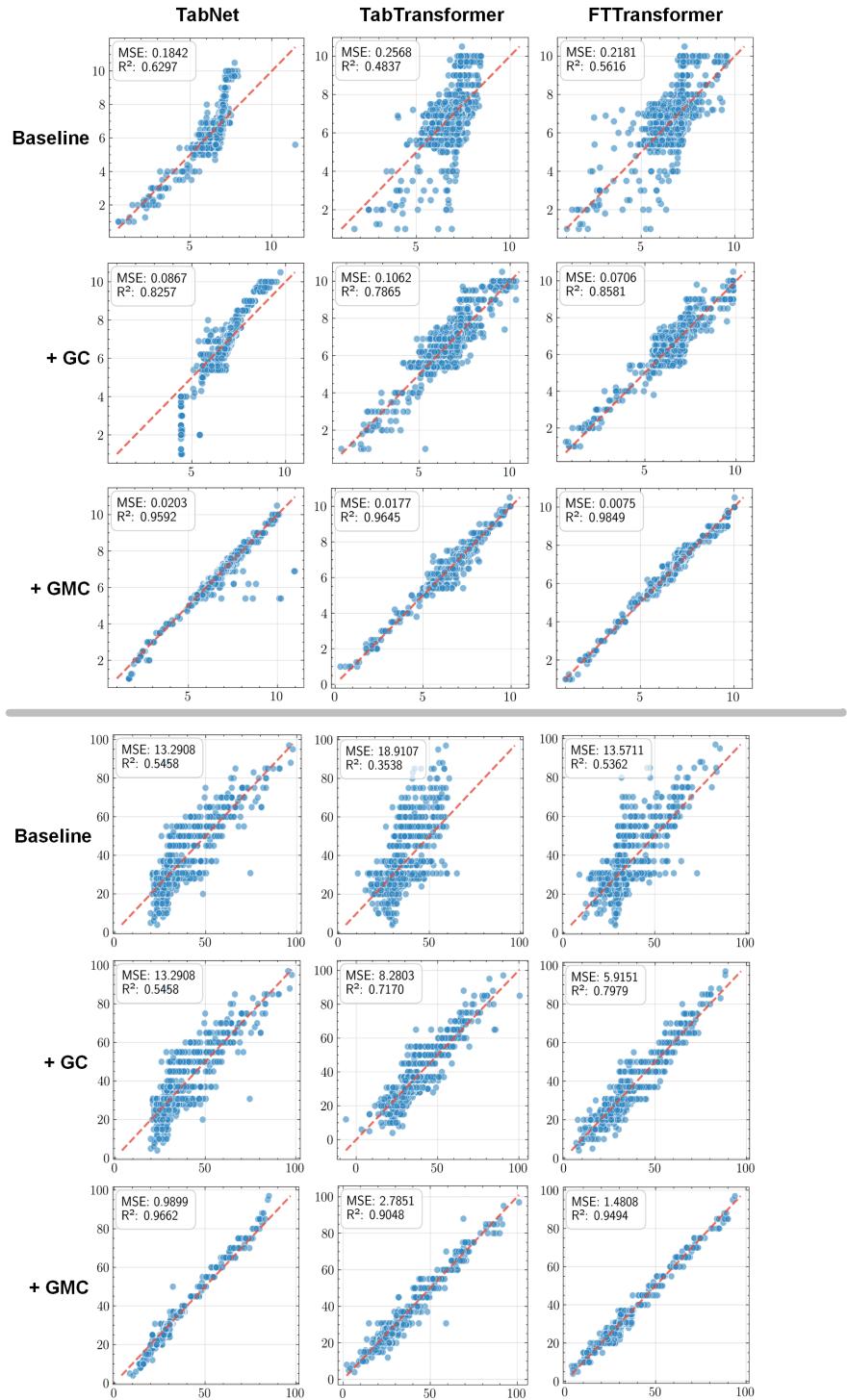


Figure 13: Illustration of the regression performance improvement example in TabNet, TabTransformer and FT-Transformer by adding GC or GMC. Each plots predicted vs. actual value.

- 513           • The authors should reflect on the factors that influence the performance of the approach.  
 514           For example, a facial recognition algorithm may perform poorly when image resolution  
 515           is low or images are taken in low lighting. Or a speech-to-text system might not be  
 516           used reliably to provide closed captions for online lectures because it fails to handle  
 517           technical jargon.
- 518           • The authors should discuss the computational efficiency of the proposed algorithms  
 519           and how they scale with dataset size.
- 520           • If applicable, the authors should discuss possible limitations of their approach to  
 521           address problems of privacy and fairness.
- 522           • While the authors might fear that complete honesty about limitations might be used by  
 523           reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
 524           limitations that aren't acknowledged in the paper. The authors should use their best  
 525           judgment and recognize that individual actions in favor of transparency play an impor-  
 526           tant role in developing norms that preserve the integrity of the community. Reviewers  
 527           will be specifically instructed to not penalize honesty concerning limitations.

### 528           3. Theory assumptions and proofs

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

531           Answer: [NA]

532           Justification: Our research presents a practical approach to mixing configurations for down-  
 533           stream predictions. No novel theoretical claims are made that require formal proof.

534           Guidelines:

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

### 545           4. Experimental result reproducibility

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

549           Answer: [Yes]

550           Justification: To ensure complete reproducibility, we provide all necessary information  
 551           in Section 3, Section 4 and Appendix. It includes methodologies, experiment setups,  
 552           computing environment, parameter settings and other implementation details, enabling  
 553           independent verification of all our claims and conclusions.

554           Guidelines:

- 555           • The answer NA means that the paper does not include experiments.
- 556           • If the paper includes experiments, a No answer to this question will not be perceived  
 557           well by the reviewers: Making the paper reproducible is important, regardless of  
 558           whether the code and data are provided or not.
- 559           • If the contribution is a dataset and/or model, the authors should describe the steps taken  
 560           to make their results reproducible or verifiable.
- 561           • Depending on the contribution, reproducibility can be accomplished in various ways.  
 562           For example, if the contribution is a novel architecture, describing the architecture fully  
 563           might suffice, or if the contribution is a specific model and empirical evaluation, it may  
 564           be necessary to either make it possible for others to replicate the model with the same  
 565           dataset, or provide access to the model. In general, releasing code and data is often

566 one good way to accomplish this, but reproducibility can also be provided via detailed  
567 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
568 of a large language model), releasing of a model checkpoint, or other means that are  
569 appropriate to the research performed.

- 570 • While NeurIPS does not require releasing code, the conference does require all submissions  
571 to provide some reasonable avenue for reproducibility, which may depend on the  
572 nature of the contribution. For example
- 573 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
574 to reproduce that algorithm.
- 575 (b) If the contribution is primarily a new model architecture, the paper should describe  
576 the architecture clearly and fully.
- 577 (c) If the contribution is a new model (e.g., a large language model), then there should  
578 either be a way to access this model for reproducing the results or a way to reproduce  
579 the model (e.g., with an open-source dataset or instructions for how to construct  
580 the dataset).
- 581 (d) We recognize that reproducibility may be tricky in some cases, in which case  
582 authors are welcome to describe the particular way they provide for reproducibility.  
583 In the case of closed-source models, it may be that access to the model is limited in  
584 some way (e.g., to registered users), but it should be possible for other researchers  
585 to have some path to reproducing or verifying the results.

## 586 5. Open access to data and code

587 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
588 tions to faithfully reproduce the main experimental results, as described in supplemental  
589 material?

590 Answer: [Yes]

591 Justification: We have made our code and data publicly accessible through the GitHub links  
592 provided in this paper.

593 Guidelines:

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

## 613 6. Experimental setting/details

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

617 Answer: [Yes]

618 Justification: We specify all the Implementation details and experimental setup in Sec-  
619 tion 4.1.

620 Guidelines:

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

626 **7. Experiment statistical significance**

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

629 Answer: [Yes]

630 Justification: The paper includes error bars for key results (e.g., Table 1), clearly stating they  
631 represent standard deviation over multiple runs with different seeds.

632 Guidelines:

- 633 • The answer NA means that the paper does not include experiments.  
634 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
635 dence intervals, or statistical significance tests, at least for the experiments that support  
636 the main claims of the paper.  
637 • The factors of variability that the error bars are capturing should be clearly stated (for  
638 example, train/test split, initialization, random drawing of some parameter, or overall  
639 run with given experimental conditions).  
640 • The method for calculating the error bars should be explained (closed form formula,  
641 call to a library function, bootstrap, etc.)  
642 • The assumptions made should be given (e.g., Normally distributed errors).  
643 • It should be clear whether the error bar is the standard deviation or the standard error  
644 of the mean.  
645 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
646 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
647 of Normality of errors is not verified.  
648 • For asymmetric distributions, the authors should be careful not to show in tables or  
649 figures symmetric error bars that would yield results that are out of range (e.g. negative  
650 error rates).  
651 • If error bars are reported in tables or plots, The authors should explain in the text how  
652 they were calculated and reference the corresponding figures or tables in the text.

653 **8. Experiments compute resources**

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

657 Answer: [Yes]

658 Justification: We detail the experimental environment in Section 4.1 and compare the time  
659 of execution between different clustering methods in Fig. 11.

660 Guidelines:

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

669 **9. Code of ethics**

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

672                  Answer: [Yes]

673                  Justification: Our research conforms with every aspect of the NeurIPS Code of Ethics.

674                  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).

## 680                  10. Broader impacts

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

683                  Answer: [NA]

684                  Justification: Our research primarily contributes to improving technical aspects of down-  
685                  stream prediction tasks and does not have broader societal implications.

686                  Guidelines:

- 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).

## 709                  11. Safeguards

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

713                  Answer: [NA]

714                  Justification: The models and data presented in our work do not pose any risks of misuse.

715                  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.

- 723           • We recognize that providing effective safeguards is challenging, and many papers do  
724           not require this, but we encourage authors to take this into account and make a best  
725           faith effort.

726           **12. Licenses for existing assets**

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

730           Answer: [Yes]

731           Justification: For every dataset used in our research, we cite its original papers or official  
732           websites. We properly credit all open-source packages used (*e.g.* pytorch).

733           Guidelines:

- 734           • The answer NA means that the paper does not use existing assets.  
735           • The authors should cite the original paper that produced the code package or dataset.  
736           • The authors should state which version of the asset is used and, if possible, include a  
737           URL.  
738           • The name of the license (*e.g.*, CC-BY 4.0) should be included for each asset.  
739           • For scraped data from a particular source (*e.g.*, website), the copyright and terms of  
740           service of that source should be provided.  
741           • If assets are released, the license, copyright information, and terms of use in the  
742           package should be provided. For popular datasets, [paperswithcode .com/datasets](https://paperswithcode.com/datasets)  
743           has curated licenses for some datasets. Their licensing guide can help determine the  
744           license of a dataset.  
745           • For existing datasets that are re-packaged, both the original license and the license of  
746           the derived asset (if it has changed) should be provided.  
747           • If this information is not available online, the authors are encouraged to reach out to  
748           the asset's creators.

749           **13. New assets**

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

752           Answer: [Yes]

753           Justification: We have included README files in our released code repositories to provide  
754           clear and comprehensive documentation.

755           Guidelines:

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

764           **14. Crowdsourcing and research with human subjects**

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

768           Answer: [NA]

769           Justification: Our work does not involve crowdsourcing or research with human subjects.

770           Guidelines:

- 771           • The answer NA means that the paper does not involve crowdsourcing nor research with  
772           human subjects.

- 773           • Including this information in the supplemental material is fine, but if the main contribu-  
774           tion of the paper involves human subjects, then as much detail as possible should be  
775           included in the main paper.  
776           • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
777           or other labor should be paid at least the minimum wage in the country of the data  
778           collector.

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

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

785           Answer: [NA]

786           Justification: Our work does not involve human experiments or study participants.

787           Guidelines:

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

798           **16. Declaration of LLM usage**

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

803           Answer: [NA]

804           Justification: LLMs are not involved in core method development of our research.

805           Guidelines:

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