# MIXING CONFIGURATIONS FOR DOWNSTREAM PREDICTION

**Anonymous authors** 

000

001

002 003 004

006

008

010 011

012

013

014

015

016

017

018

019

021

023

025 026 027

028 029

031

033

034

037

040

041

042

043

044

046

047

048

049

051

052

Paper under double-blind review

#### **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, vet 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<sup>2</sup> 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.

#### 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 (MacQueen, 1967; Jianbo Shi & Malik, 2000; Ng et al., 2001). 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 (Quinn & Eimas, 1996; Bornstein et al., 2010). In particular, they often organize them at different abstraction levels (Zaadnoordijk et al., 2022; Muttenthaler et al., 2024). Inspired by this, recent advances in community detection have extended clustering to the discovery of configurations—hierarchical clusterings that span multiple resolution scales (Pitsianis et al., 2023). For example, as illustrated in the lineage diagram of Fig. 1, in the CIFAR10 dataset (Krizhevsky, 2009), coarse configurations may separate vehicles from

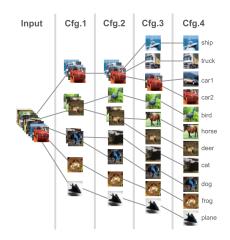


Figure 1: Illustration of CIFAR10 configurations. Each column represents a configuration—clustering at a specific resolution.

animals, while finer configurations distinguish between birds, cats, and dogs. These multi-resolution representations reveal rich hierarchical structures that could provide stronger priors or inductive biases for deep models. However, despite their potential, such configurations remain largely underexplored in deep learning, especially in challenging domains where labels are sparse.

One such domain is 16S ribosomal RNA (rRNA) gene sequencing, a widely used tool in microbiome studies for identifying and classifying bacteria. Analyzing 16S rRNA data has consistently confronted significant challenges in downstream prediction tasks within label-scarce environments.

Previous works in 16S rRNA representation learning have demonstrated substantial benefits for bacterial taxonomic profiling and microbial community analysis (Janda & Abbott, 2007; Wang et al., 2007; De Vrieze et al., 2018). Notably, Johnson et al. (2019) showed that full-length sequencing combined with appropriate clustering of intragenomic sequence variation can provide more accurate representation of bacterial species in microbiome datasets. These findings underscore the importance of learning clustered representations without relying on labels.

Recent methodologies typically transform clustering results into pseudo-labels to enhance downstream prediction performance. For instance, DeepCluster (Caron et al., 2019) iteratively clusters CNN-extracted visual features and leverages these cluster assignments to guide network parameter updates. Graph-based methods such as (Yang et al., 2023) employ structural clustering to overcome limitations of traditional contrastive learning approaches that depend on positive and negative sample pairs. Their method captures structural relationships among nodes in heterogeneous information networks, establishing a self-supervised pre-training framework that learns robust network representations from unlabeled data. Nevertheless, these approaches predominantly focus on a single configuration, overlooking the potential benefits of mixing configurations across multiple resolution scales.

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

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

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

#### 2 Preliminary results

We first present preliminary experimental results on configurations using CIFAR10. Specifically, we compare patterns of configurations with those of the learnable "register" tokens in a recent vision transformer DINOv2-reg (Darcet et al., 2024). Fig. 2 shows the attention maps from our configurations and their register tokens. Moreover, Fig. 3 shows qualitative behaviors of our configurations and their quantitative advantages over registers in terms of feature importance and neighborhood similarity. From these results, we identify three key properties:

Configurations emerge via unsupervised or self-supervised learning. We define Near ground truth (GT) balls as balls selected with the highest clustering scores, marked yellow in Fig. 2a. As shown in Fig. 2b, the attention map, acquired by feeding configurations as tokens to attention heads for linear probing, yields high norm regions substantially overlap with GT balls. On another hand, DINOv2-reg exhibits similar attention map patterns in selected registers (see Fig. 2c), which might be related to registers activating different areas in Fig. 2d, similar to slot attention (Locatello et al., 2020; Caron et al., 2021; Oquab et al., 2024; Darcet et al., 2024). Thus, based on the similar attention map behavior, register token can be considered as a latent configuration.

Configurations are selected and mixed based on input and task. Configuration selection and mixing refers to learning which resolution scales to focus on for a given downstream task. We visualize this via attention maps over configuration tokens, where high-norm regions indicate the selected scales. In Fig. 2b, attention norms vary across rows, showing that each input sample triggers different resolution scales. Without any change to the configurations, we merge the original labels into coarser classes (Fig. 3a) and plot the new attention map (Fig. 3b). The attention shifts to align with the coarser GT, whereas DINOv2-reg register tokens remain unchanged unless re-trained. These observations confirm that configuration selection and mixing are input- and task-dependent.

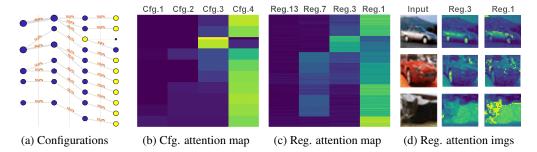


Figure 2: Comparison of attention maps obtained from configurations and registers, rows for samples. (a): Lineage diagram for configurations, near GT balls are marked yellow. (b): Attention map of configuration tokens in an attention-based linear probing. (c): Attention map of DINOv2-reg register tokens,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

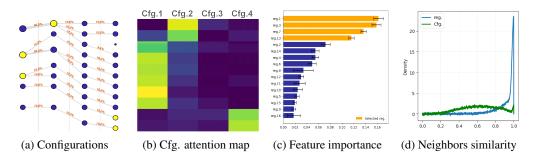


Figure 3: Illustration of another two properties of configurations, grouped by left two and right two. (a): Lineage diagram where coarser classes are used for GT. (b): Attention map in linear probing the coarser classes. (c): Distribution of feature vector importance over the register tokens querying, 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.

Configurations are more informative and less redundant than register tokens. Register tokens can help extract configurations, similar to object detection (Siméoni et al., 2021; Zhang et al., 2022), but selecting a fixed number by feature importance is arbitrary and non-rigorous (see Fig. 3c). Furthermore, register tokens exhibit high redundancy—cosine similarity between their embeddings and their 2 neighbors embeddings is heavily skewed toward 1—whereas configurations yield information less redundant (see Fig. 3d).

#### 3 METHODOLOGY

Having these characterizations, we hypothesize that unsupervised methods can produce hierarchical *multi-resolution clusterings*, and that task- and input-specific *selection and mixing* of these configurations represent *global information* beneficial to downstream tasks. Building on the hypothesis, we propose a lightweight module *GraMixC*, that treats configurations as tokens ([CFG]) and incorporates a novel alignment layer plus learnable attention heads (Vaswani et al., 2017) after the configuration extraction model, enabling task- and input-specific mixing of configurations via end-to-end back-propagation.

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

Except for the downstream predictor, the GraMixC model can be divided into three parts: unsupervised configuration learning, the Reverse Merge & Split (RMS) for alignment, and attention heads for fusion. In the attention heads part, following Darcet et al. (2024), we append register tokens [REG]

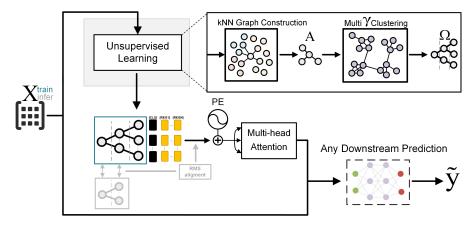


Figure 4: Illustration of the proposed GraMixC module and resulting model. The input data branches into (upper) a path to unsupervised learning box that extracts configurations, and (lower) a direct path to the downstream predictor. Their outcomes concatenate and pass to the downstream predictor. The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

after [CFG] and [CLS] for a clean attention map, that can be used backwards to guide configuration selection. Below we detail the rest two components in Section 3.1 and Section 3.2.

#### 3.1 MULTI-RESOLUTION GRAPH-BASED CLUSTERING

Given X, multi-resolution clustering seeks to extract configurations—valid hierarchical clusterings across multiple resolution scales—which we denote as  $\Omega \in \mathbb{N}^{N \times m}$ , where m denotes the number of valid resolution levels. To preserve the latent manifold structure in data, ease parameter sensitivity, and prevent other problems with traditional clustering methods (see Section D), we choose the resolution parameter ( $\gamma \in \mathbb{R}_+$ )-based community detection as our core clustering method. While BlueRed (Liu et al., 2021) can conduct graph clustering without problems like resolution limit or parameter sensitivity in traditional methods, recent work by Pitsianis et al. (2023) further demonstrates the elimination of  $\gamma$  selection, and enabled the unsupervised discovery of  $\Omega$  and the corresponding set of all valid  $\gamma$ , which is denoted as  $\Gamma = \{\gamma_1^*, \gamma_2^*, \ldots, \gamma_m^*\} \subseteq [0, \infty)$ . Inspired by these works, the unsupervised box in Fig. 4 unfolds into two steps: (1) k-nearest neighbors (kNN) (Tenenbaum et al., 2000) graph construction, which return a directed graph G = (V, E), usually represented as adjacency matrix  $A \in \mathbb{R}_+^{N \times N}$ , and (2) multi- $\gamma$  clustering on the resulted graph, *i.e.* modularity based community detection with unsupervised  $\Gamma$  learning, which return the wanted  $\Omega$ . The details for each of these two steps are:

(1) kNN graph construction. We construct a kNN graph with  $k = \log_{10} N$  as convention, using Euclidean distance for simplicity. Such pair-wise geometric distance between two different vertexes is denoted  $d(\boldsymbol{x}_i, \boldsymbol{x}_j)$  where  $i \neq j$  and  $x_i \in \mathbb{R}^d$  is the i-th feature vector. We then have the adjacency matrix  $\boldsymbol{A}$  formulated as:  $A_{ij} = d(\boldsymbol{x}_i, \boldsymbol{x}_j)$  if  $(\boldsymbol{x}_i, \boldsymbol{x}_j) \in E$ , 0 otherwise, where E is the edge set of the kNN graph and  $A_{ij}$  denotes the i-th row and j-th column element of the adjacency matrix. Then we force column stochastic by dividing each column in the constructed  $\boldsymbol{A}$  with the column sum. The resulted graph is sparse stochastic, and we can apply Stochastic Graph t-SNE (SG-t-SNE) reweighting (Pitsianis et al., 2019), which proved to remedy skewed degree distribution, that is not promised by conventional t-SNE (Van der Maaten & Hinton, 2008). From the original work, the key equations for SG-t-SNE reweighting are:

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

where  $\lambda$  is a non-negative constant, which we simply set to 15 as previous work show that it is not so sensitive to the choice of  $\lambda$  (Pitsianis et al., 2019), and  $\sigma_i$  is a variable to be numerically solved with bisection method. After giving value of w to d, we have A with less skewed degree distribution, which avoids problems like numerical instability and bias towards hubs in downstream clustering.

(2) multi- $\gamma$  community detection. Then one may simply pass the reweighted A to  $\gamma$ -based community detection method, such as Leiden algorithm (Traag et al., 2019), to get one pseudo-configuration vector  $\boldsymbol{\omega}_{\gamma} \in \{1,\dots,N\}^N$  ("pseudo" for not sure to be valid). However, such  $\gamma$  falls in the range of  $[0,\infty)$ , and searching over all possible  $\gamma$  is exhausting. Therefore, we incorporate the BlueRed method with parallel descending triangulation (parallel-DT) (Pitsianis et al., 2023), in order to automatically discover all valid  $\gamma^* \in \Gamma$ . Given a fixed  $\gamma$ , BlueRed find the optimal configuration  $\boldsymbol{\omega}_{\gamma}$  by the following optimization:

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

where  $\omega_i$  denotes the i-th element of  $\omega$ ,  $|\omega|_{\infty} = \max_{i \leq N} \omega_i$  is a inf-norm, and 1 denotes the indicator gate which take value 1 if its subscript condition holds, 0 otherwise. Pitsianis et al. (2023) describe the first term as attraction and the second term as repulsion. Optimizing each solely yields all-in-one configuration  $\omega_0 = [1, 1, \dots, 1]$  and all-lonely configuration  $\omega_{\infty} = [1, 2, \dots, N]$ . Between these two configurations, parallel-DT allows forming BlueRed Front (BRF) (Pitsianis et al., 2023) by segmenting  $(0, \infty)$  into m ranges, among which each has a dominant  $\gamma_i^*$  yields lower HAR (Pitsianis et al., 2023)—the sum of first term and the negative second term—which means "local minimum" on that range. Thus desired  $\Omega$  is formed.

#### 3.2 RMS: REVERSE MERGE & SPLIT ALIGNMENT

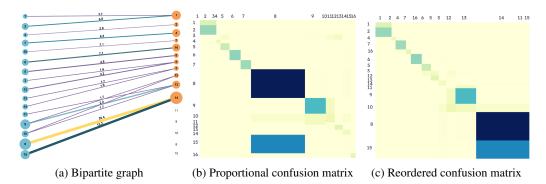


Figure 5: Example of the RMS alignment process applied to clustering results and ground truth (both treated as configurations) on the Salinas dataset (Plaza & Tilton, 2005). (a): Bipartite graph representation, where blue nodes correspond to predicted clusters and red nodes to ground truth clusters. 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  $C_{\rm tw}$  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.

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

$$SCORE(\boldsymbol{\omega}_i, \boldsymbol{\omega}_j) = ARI(\boldsymbol{\omega}_i, \boldsymbol{\omega}_j) - \theta \left| \frac{|\boldsymbol{\omega}_i|_{\infty} - |\boldsymbol{\omega}_j|_{\infty}}{|\boldsymbol{\omega}_i|_{\infty} + |\boldsymbol{\omega}_j|_{\infty}} \right|.$$

where  $\theta$  is a hyperparameter to balance the weights of the two terms, which we set to 0.1, ARI is the adjusted rand index as defined in Hubert & Arabie (1985). By this punished ARI design, we consider different labels, merge and split during scoring the alignment between two partition, but also avoids too much difference in number of clusters (one extreme case is  $\omega_0$  and  $\omega_\infty$  has ARI of 1).

However, the SCORE itself does not convey the mapping we need for reassigning. In *RMS* alignment, 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 process with a concrete example, showing how the confusion matrix captures the relationship between predicted and ground truth clusters, including cases where clusters are split or merged across configurations. As an assignment problem with a rectangle cost matrix  $-C^1$ , it is solvable by twisting existing Hungarian algorithm methods (Kuhn, 1955; Jonker & Volgenant, 1987; Bertsekas, 1992). 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 (Fiedler, 1973) assumes symmetric positive semi-definite, it is not directly applicable to C. Inspired by a recent work of Floros et al. (2024), we introduce a *two-walk Laplacian*, which is defined as:

$$m{L}_{ ext{tw}} = m{D} - m{C}_{ ext{tw}}, \quad ext{with} \quad m{C}_{ ext{tw}} = \left[ egin{array}{cc} m{C}m{C}^ op & m{C} \\ m{C}^ op & m{C}^ op m{C} \end{array} 
ight],$$

where  $D = \operatorname{diag}(C_{\operatorname{tw}} 1)$  is the diagonal degree matrix of  $C_{\operatorname{tw}}$ . 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  $L_{\operatorname{tw}}$ , 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  $\Omega_{\text{train}}$  and  $\Omega_{\text{test}}$  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 our GitHub repository  $^2$  and Section E 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

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 (Kingma & Ba, 2017) 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 (German Collection of Microorganisms and Cell Cultures GmbH, 2025) and NIH. 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. (2024) and related works (Wood & Salzberg, 2014; Compeau et al., 2011), we encode each 16S rRNA sequence as a 7-mer count vector in  $\mathbb{N}^{16\,384}$ , yielding a dataset of shape  $65\,023\times16\,384$ . 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. Section C provides an illustration for target value ( $y_{\text{train}}$  and  $y_{\text{test}}$ ) distribution. Preprocessing—robust scaling, variance thresholding, and selection

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

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

Table 1: Regression performance on DSNI-pH, DSNI-Temp and QM9. Values are mean±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 ↓               | $R^2$               | MSE ↓               | $R^2$               | MAE ↓               | $\mathbb{R}^2$      |
| RF       | 0.198±0.000         | 0.601±0.001         | 17.759±0.276        | 0.393±0.009         | 0.015±0.000         | 0.979±0.000         |
| XGBoost  | 0.196±0.001         | 0.604±0.003         | 18.212±0.543        | 0.377±0.018         | 0.014±0.001         | 0.978±0.001         |
| CatBoost | 0.193±0.001         | 0.610±0.002         | 17.375±0.398        | 0.406±0.013         | 0.014±0.000         | 0.978±0.002         |
| 3LP      | 0.201±0.002         | 0.595±0.006         | 18.484±0.183        | 0.368±0.006         | 0.018±0.001         | 0.958±0.001         |
| 3LP+GC   | 0.097±0.004         | 0.804±0.008         | 6.520±0.360         | 0.777±0.012         | 0.016±0.003         | 0.974±0.000         |
| 3LP+GMC  | 0.023±0.002         | <b>0.953</b> ±0.004 | 2.277±0.061         | <b>0.922</b> ±0.002 | <b>0.010</b> ±0.003 | <b>0.990</b> ±0.002 |
| TabN     | 0.184±0.004         | 0.629±0.007         | 13.290±0.244        | 0.545±0.008         | 0.015±0.001         | 0.962±0.002         |
| TabN+GC  | 0.086±0.003         | 0.825±0.007         | 7.997±0.210         | 0.726±0.007         | 0.012±0.002         | 0.983±0.001         |
| TabN+GMC | <b>0.020</b> ±0.001 | <b>0.959</b> ±0.002 | <b>0.989</b> ±0.361 | <b>0.966</b> ±0.012 | <b>0.008</b> ±0.000 | <b>0.995</b> ±0.002 |
| TabT     | 0.256±0.007         | 0.483±0.014         | 18.910±0.247        | 0.353±0.008         | 0.434±0.008         | 0.921±0.008         |
| TabT+GC  | 0.106±0.002         | 0.786±0.005         | 8.280±0.303         | 0.717±0.010         | 0.212±0.004         | 0.961±0.008         |
| TabT+GMC | <b>0.017</b> ±0.002 | <b>0.964</b> ±0.005 | <b>2.785</b> ±0.540 | <b>0.904</b> ±0.018 | <b>0.009</b> ±0.000 | <b>0.998</b> ±0.001 |
| FTT      | 0.218±0.003         | 0.561±0.006         | 13.571±0.069        | 0.536±0.002         | 0.085±0.005         | 0.984±0.006         |
| FTT+GC   | 0.070±0.003         | 0.858±0.007         | 5.915±0.277         | 0.797±0.009         | 0.034±0.002         | 0.993±0.003         |
| FTT+GMC  | <u>0.007</u> ±0.005 | <b>0.984</b> ±0.009 | <b>1.480</b> ±0.120 | <b>0.949</b> ±0.004 | <b>0.026</b> ±0.001 | <b>0.995</b> ±0.003 |

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 (Ramakrishnan et al., 2014) for molecular property regression, on Boston Housing (Harrison & Rubinfeld, 1978), and on MNIST (Lecun et al., 1998) and CIFAR10 for classification (some in Section F).

**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) (Breiman, 2001), XGBoost (Chen & Guestrin, 2016), CatBoost (Prokhorenkova et al., 2018). 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 hidden dims [256,128,64]. Because our inputs combine numerical features with categorical configurations, we naturally consider tabular models: TabNet (TabN) (Arik & Pfister, 2020), TabTransformer (TabT) (Huang et al., 2020), FT-Transformer (FTT) (Gorishniy et al., 2023) were all run with their default settings from the official implementations.

#### 4.2 EVALUATION OF THE PROPOSED MODULE

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

**Performance improvement.** Table 1 shows that adding GC and GMC yields consistent gains across all baselines. Among these observed improvements, the scores increasing on DSNI is quite satisfying. Prior specialized growth-media regression methods are not convincing with  $R^2 \leq 0.8$  (e.g., 0.75 (Sauer & Wang, 2019)). We confirm this with our base models score  $R^2$  between 0.3 and 0.6 on DSNI-pH and DSNI-Temp. However, even without tailoring the baseline model design, we bring the score to a new high by simply adding GC or GMC. Fig. 6 illustrates some examples

of such improvement. We see the model's predictions align more closely with the ideal regression line and better handle rare cases, by incorporating configurations and probably capturing the latent manifold structure. Incorporating GC and further GMC raises  $R^2$  to 0.98 (pH) and 0.97 (Temp). Which not only is considered very satisfying in application of bacterial cultivation but also set the new state-of-the-art (SOTA) for growth-media prediction. On QM9, GraMixC achieves an MAE of 0.008, nearly matching the SOTA (w/o extra training data) of 0.007 (Fang et al., 2022), and represents the best result among non-GNN models.

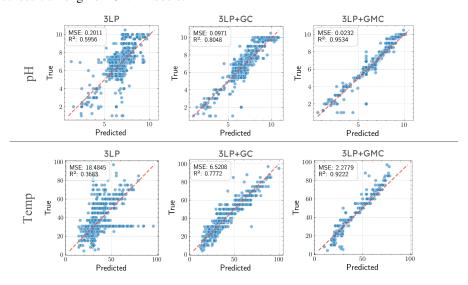


Figure 6: Illustration of the regression performance improvement example in 3LP by adding GC or GMC. Each column plots predicted vs. actual pH (top) or temperature (bottom). 3LP+GC (middle) outperforms the 3LP baseline (left), while 3LP+GMC (right) further boosts  $R^2$  up to > 0.9.

**Number of configurations used.** We ablate the number of configuration levels in GMC. Fig. 7 shows that more configurations generally decreases MSE and increases R<sup>2</sup>, confirming the value of multi-resolution information. Importantly, GMC often needs more than half as many total configurations to outperform GC, and performance plateaus—or even slightly declines—when including the last few configurations. These aligns with Pitsianis et al. (2023), who report a finite set of optimal configurations rather than continuous gains at infinite resolutions. Using all configurations available is still preferred.

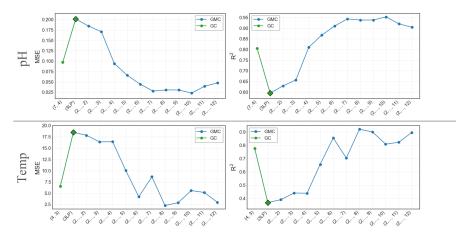
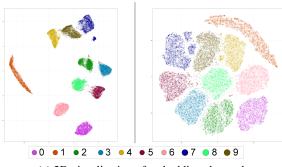


Figure 7: Ablation study on the number of configurations used on DSNI. On the blue curves (GMC),  $[2,\ldots,i]$  denote fusing configurations from 2 through i via GraMixC. On the green curves (GC), (i,j) denote the best train/test configuration pair used in static concatenation. Incrementally mixing configurations improves performance and outperforms static concatenation.

### 4.3 QUALITATIVE EVALUATION OF CONFIGURATIONS.

Our final experiment compares configurations against standard representation-extraction methods. As discussed in Section 1, configurations can be viewed as special unsupervised representation learning. Fig. 3 already shows their advantage over self-supervised register tokens. Here, we replace GC/GMC with PCA (Jolliffe & Cadima, 2016), UMAP (McInnes et al., 2020), and a vanilla autoencoder (AE), each embed into dimensions the same number of as our configurations. We visualize these embeddings on MNIST (Fig. 8a; additional views in Section F.2). Qualitatively, SG-t-SNE (the reduction step in GraMixC) yields more uniform, well-separated clusters that respect global kNN connectivity rather than forming hubs. Fig. 8b quantifies downstream classification accuracy, where GC and GMC strongly outperform PCA, UMAP, and AE given the same embedding budget. These results confirm that mixed configurations provide a more expressive yet compact representation for downstream tasks.



|          | $CE\downarrow$ | 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 |
| 3LP+GMC  | 0.028          | 0.993 |

(a) 2D visualization of embeddings learned.

(b) Classification performance.

Figure 8: (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.

### 5 CONCLUSION

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

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

# ETHICS STATEMENT

We followed the ICLR Code of Ethics. All datasets (DSNI, NIH/DSMZ metadata, QM9, Boston Housing, MNIST, CIFAR-10) are publicly available and contain no personally identifiable information. Use complies with their licenses.

**Risks:** Our method predicts cultivation conditions from 16S rRNA features. Outputs are not lab-ready instructions and require expert validation. We discourage unsafe or unsupervised use, particularly with pathogenic organisms.

**Bias:** Training data reflect known biases (e.g., over-represented mesophiles). We report distributions (Section C) and evaluate across multiple benchmarks to reduce overfitting.

**Conflicts:** No competing financial interests. Experiments were run on institutional hardware.

#### REPRODUCIBILITY STATEMENT

Code, configs, and data-processing scripts are available at project-34CB and project-82CE on https://anonymous.4open.science/r/project-34CB and https://anonymous.4open.science/r/project-82CE.

All algorithms, hyperparameters, and dataset details are given in Sections C, 3.1, 3.2 and 4.1. Splits and seeds are fixed and provided. Figures and tables can be reproduced directly from the released scripts. Compute setup (GPU, runtime, nondeterminism) is documented in the README.

#### REFERENCES

- Sercan O. Arik and Tomas Pfister. TabNet: Attentive Interpretable Tabular Learning, December 2020.
- Dimitri P. Bertsekas. Auction algorithms for network flow problems: A tutorial introduction. *Computational Optimization and Applications*, 1(1):7–66, October 1992. ISSN 0926-6003, 1573-2894. doi: 10.1007/BF00247653.
- Marc H. Bornstein, Martha E. Arterberry, and Clay Mash. Infant object categorization transcends diverse object-context relations. *Infant Behavior & Development*, 33(1):7–15, February 2010. ISSN 1934-8800. doi: 10.1016/j.infbeh.2009.10.003.
- Leo Breiman. Random Forests. *Machine Learning*, 45(1):5–32, 2001. ISSN 1573-0565. doi: 10.1023/A:1010933404324.
- Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep Clustering for Unsupervised Learning of Visual Features, March 2019.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers, May 2021.
- Abdulkadir Çelikkanat, Andres R Masegosa, and Thomas D Nielsen. Revisiting K-mer profile for effective and scalable genome representation learning, November 2024.
- Tianqi Chen and Carlos Guestrin. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, August 2016. doi: 10.1145/2939672.2939785.
- Phillip E. C. Compeau, Pavel A. Pevzner, and Glenn Tesler. How to apply de Bruijn graphs to genome assembly. *Nature Biotechnology*, 29(11):987–991, November 2011. ISSN 1546-1696. doi: 10.1038/nbt.2023.
- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers, April 2024.
- Jo De Vrieze, Ameet J. Pinto, William T. Sloan, and Umer Zeeshan Ijaz. The active microbial community more accurately reflects the anaerobic digestion process: 16S rRNA (gene) sequencing as a predictive tool. *Microbiome*, 6(1):63, April 2018. ISSN 2049-2618. doi: 10.1186/s40168-018-0449-9.

Xiaomin Fang, Lihang Liu, Jieqiong Lei, Donglong He, Shanzhuo Zhang, Jingbo Zhou, Fan Wang, Hua Wu, and Haifeng Wang. Geometry-enhanced molecular representation learning for property prediction. *Nature Machine Intelligence*, 4(2):127–134, February 2022. ISSN 2522-5839. doi: 10.1038/s42256-021-00438-4.

Miroslav Fiedler. Algebraic connectivity of graphs. *Czechoslovak Mathematical Journal*, 23(2): 298–305, 1973. ISSN 0011-4642, 1572-9141. doi: 10.21136/CMJ.1973.101168.

- Dimitris Floros, Nikos Pitsianis, and Xiaobai Sun. Algebraic vertex ordering of a sparse graph for adjacency access locality and graph compression. In 2024 IEEE High Performance Extreme Computing Conference (HPEC), pp. 1–7, September 2024. doi: 10.1109/HPEC62836.2024. 10938496.
- German Collection of Microorganisms and Cell Cultures GmbH. DSMZ german collection of microorganisms and cell cultures, 2025.
- Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting Deep Learning Models for Tabular Data, October 2023.
- David Harrison and Daniel L Rubinfeld. Hedonic housing prices and the demand for clean air. *Journal of Environmental Economics and Management*, 5(1):81–102, March 1978. ISSN 0095-0696. doi: 10.1016/0095-0696(78)90006-2.
- Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. TabTransformer: Tabular data modeling using contextual embeddings, December 2020.
- Lawrence J. Hubert and Phipps Arabie. Comparing partitions. *Journal of Classification*, 2(2-3): 193–218, 1985. ISSN 0176-4268. doi: 10.1007/BF01908075.
- J. Michael Janda and Sharon L. Abbott. 16S rRNA gene sequencing for bacterial identification in the diagnostic laboratory: Pluses, perils, and pitfalls. *Journal of Clinical Microbiology*, 45(9): 2761–2764, September 2007. ISSN 0095-1137. doi: 10.1128/JCM.01228-07.
- Jianbo Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, August 2000. ISSN 01628828. doi: 10.1109/34.868688.
- Jethro S. Johnson, Daniel J. Spakowicz, Bo-Young Hong, Lauren M. Petersen, Patrick Demkowicz, Lei Chen, Shana R. Leopold, Blake M. Hanson, Hanako O. Agresta, Mark Gerstein, Erica Sodergren, and George M. Weinstock. Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis. *Nature Communications*, 10(1):5029, November 2019. ISSN 2041-1723. doi: 10.1038/s41467-019-13036-1.
- Ian T. Jolliffe and Jorge Cadima. Principal component analysis: A review and recent developments. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, April 2016. doi: 10.1098/rsta.2015.0202.
- R. Jonker and A. Volgenant. A shortest augmenting path algorithm for dense and sparse linear assignment problems. *Computing*, 38(4):325–340, December 1987. ISSN 1436-5057. doi: 10.1007/BF02278710.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, January 2017.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Toronto, ON, Canada, 2009.
- H. W. Kuhn. The Hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 2(1-2):83–97, March 1955. ISSN 0028-1441, 1931-9193. doi: 10.1002/nav.3800020109.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998. ISSN 1558-2256. doi: 10.1109/5.726791.

- Tiancheng Liu, Dimitris Floros, Nikos Pitsianis, and Xiaobai Sun. Digraph Clustering by the BlueRed Method. In 2021 IEEE High Performance Extreme Computing Conference (HPEC), pp. 1–7, September 2021. doi: 10.1109/HPEC49654.2021.9622834.
  - Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention, October 2020.
  - J. MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, volume 5.1, pp. 281–298. University of California Press, January 1967.
  - Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction, September 2020.
  - Lukas Muttenthaler, Klaus Greff, Frieda Born, Bernhard Spitzer, Simon Kornblith, Michael C. Mozer, Klaus-Robert Müller, Thomas Unterthiner, and Andrew K. Lampinen. Aligning machine and human visual representations across abstraction levels, October 2024.
  - Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001. doi: 10.5555/2980539.2980649.
  - Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision, February 2024.
  - Nikos Pitsianis, Alexandros-Stavros Iliopoulos, Dimitris Floros, and Xiaobai Sun. Spaceland Embedding of Sparse Stochastic Graphs. In *2019 IEEE High Performance Extreme Computing Conference (HPEC)*, pp. 1–8, September 2019. doi: 10.1109/HPEC.2019.8916505.
  - Nikos Pitsianis, Dimitris Floros, Tiancheng Liu, and Xiaobai Sun. Parallel Clustering with Resolution Variation. In *2023 IEEE High Performance Extreme Computing Conference (HPEC)*, pp. 1–8, Boston, MA, USA, September 2023. IEEE. ISBN 979-8-3503-0860-0. doi: 10.1109/HPEC58863. 2023.10363552.
  - A.J. Plaza and J.C. Tilton. Automated selection of results in hierarchical segmentations of remotely sensed hyperspectral images. In *Proceedings*. 2005 IEEE International Geoscience and Remote Sensing Symposium, 2005. IGARSS '05., volume 7, pp. 4946–4949, July 2005. doi: 10.1109/IGARSS.2005.1526784.
  - Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. CatBoost: Unbiased boosting with categorical features. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
  - P. C. Quinn and P. D. Eimas. Perceptual cues that permit categorical differentiation of animal species by infants. *Journal of Experimental Child Psychology*, 63(1):189–211, October 1996. ISSN 0022-0965. doi: 10.1006/jecp.1996.0047.
  - Raghunathan Ramakrishnan, Pavlo O. Dral, Matthias Rupp, and O. Anatole von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data*, 1(1):140022, August 2014. ISSN 2052-4463. doi: 10.1038/sdata.2014.22.
  - David B Sauer and Da-Neng Wang. Predicting the optimal growth temperatures of prokaryotes using only genome derived features. *Bioinformatics*, 35(18):3224–3231, September 2019. ISSN 1367-4803. doi: 10.1093/bioinformatics/btz059.
  - Oriane Siméoni, Gilles Puy, Huy V. Vo, Simon Roburin, Spyros Gidaris, Andrei Bursuc, Patrick Pérez, Renaud Marlet, and Jean Ponce. Localizing objects with self-supervised transformers and no labels, September 2021.

Joshua B. Tenenbaum, Vin de Silva, and John C. Langford. A Global Geometric Framework for Nonlinear Dimensionality Reduction. *Science*, 290(5500):2319–2323, December 2000. doi: 10.1126/science.290.5500.2319.

V. A. Traag, L. Waltman, and N. J. Van Eck. From Louvain to Leiden: Guaranteeing well-connected communities. *Scientific Reports*, 9(1):5233, March 2019. ISSN 2045-2322. doi: 10.1038/s41598-019-41695-z.

Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11):2579–2605, 2008. ISSN 1533-7928.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

Qiong Wang, George M. Garrity, James M. Tiedje, and James R. Cole. Naive bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. *Applied and Environmental Microbiology*, 73(16):5261–5267, August 2007. ISSN 0099-2240. doi: 10.1128/AEM.00062-07.

Derrick E. Wood and Steven L. Salzberg. Kraken: Ultrafast metagenomic sequence classification using exact alignments. *Genome Biology*, 15(3):R46, March 2014. ISSN 1474-760X. doi: 10.1186/gb-2014-15-3-r46.

Yaming Yang, Ziyu Guan, Zhe Wang, Wei Zhao, Cai Xu, Weigang Lu, and Jianbin Huang. Self-supervised Heterogeneous Graph Pre-training Based on Structural Clustering, April 2023.

Lorijn Zaadnoordijk, Tarek R. Besold, and Rhodri Cusack. Lessons from infant learning for unsupervised machine learning. *Nature Machine Intelligence*, 4(6):510–520, June 2022. ISSN 2522-5839. doi: 10.1038/s42256-022-00488-2.

Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M. Ni, and Heung-Yeung Shum. DINO: DETR with improved DeNoising anchor boxes for end-to-end object detection, July 2022.

#### A LLM USAGE DISCLOSURE

We used GPT-5 only for writing support: editing text, drafting boilerplate (ethics/reproducibility), and LaTeX troubleshooting. It was not used for research ideas, experiments, or results. All technical content was created and verified by the authors. No sensitive data were shared.

#### B AN INTUITIVE EXAMPLE OF CONFIGURATION MIXING

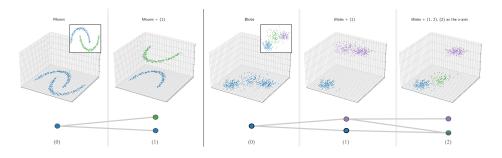


Figure 9: Illustration of multi-resolution clustering on synthetic datasets. GT is shown in the framed box in (0). Upper is the embedding of Moons (left) and Blobs (right) with corresponding configuration (i) as third dimension; lower is the lineage diagram of the configurations.

To illustrate the necessity of fusing valid clusterings across resolution scales, we use two synthetic point-cloud datasets from scikit-learn: "Moons" and "Blobs." The Blobs dataset is tuned so that no

single clustering resolution recovers all three clusters. Fig. 9 visualizes each dataset in 3D, using the third axis to encode cluster assignments for the corresponding configuration: coarser configuration (1) and finer configuration (2). Configuration (1), by lifting some dots above the plane, cleanly separates the two Moon arcs but merges two (purple and green) of the Blobs clusters. Configuration (2), by itself, fails the Blobs with a different merge (blue and green). Only by fusing both configurations can all clusters be disentangled—the purple dots in (1) that falls down in (2), emerges correct as the green cluster. This toy example shows that multi-resolution clusterings alone are insufficient without a fusion mechanism. Our GraMixC use attention-based fusion to integrate these scales. While just one demonstration, it highlights the broader advantage of mixing configurations in complex settings.

### **DSNI** DATASET DISTRIBUTION

Fig. 10 shows the pH and temperature target distributions for DSNI across training and test splits. Both splits cover similar ranges, though with natural imbalance (e.g., mesophilic temperatures, pH 6-8) reflecting biases in the underlying NIH/DSMZ data. These distributions are important for interpreting regression performance and highlight potential challenges under distributional shift.

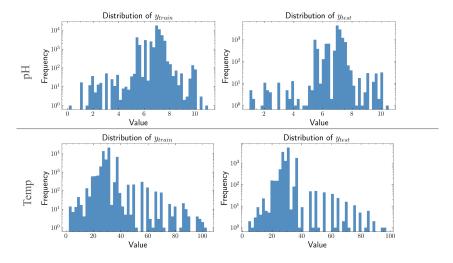


Figure 10: Illustration of target value distributions across train-test splits in DSNI dataset. The first row represents pH distributions and the second row represents temperature distributions. The first column represents the training set ( $y_{\text{train}}$ ) and the second column represents the test set ( $y_{\text{test}}$ ).

#### SYNTHETIC CLUSTERING BENCHMARKS D

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

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

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

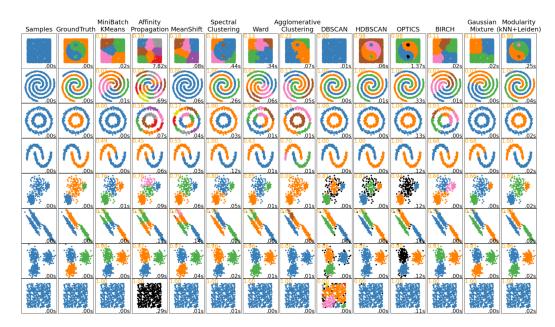


Figure 11: Illustration of clustering methods comparisons across multiple synthetic datasets. Rows correspond to different 2D point clouds—the first row is custom, others are from scikit-learn. Each method's result is labeled with ARI (top-left in yellow) and execution time (bottom-right in black). *Modularity:*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

#### E RMS ALIGNMENT DETAILS

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

### Implementation notes.

- We set  $\theta = 0.1$  and compute ARI as in Hubert & Arabie (1985).
- We use 0.1 % of the train samples as anchors to form A.
- The greedy matching loops over each train configuration  $\omega_i$  to find its best-scoring test partner  $\omega_j$ , applies the label mapping, and removes both from further consideration to ensure one-to-one alignment.

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

#### F ADDITIONAL EXPERIMENTAL RESULTS

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

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

#### F.1 ADDITIONAL EVALUATION OF PROPOSED MODULE

Table 2 extends our evaluation to three additional benchmarks: Boston Housing (regression), MNIST and CIFAR-10 (classification). We compare classical ensembles (RF, XGBoost, CatBoost), a 3-layer

841 842

843

844 845

846

847

848

849

850

851

852

853

854 855

856

858

859

861

862 863

# Algorithm 1 Reverse Merge & Split (RMS) Alignment

```
811
                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
812
                Ensure: Aligned \Omega_{\mathrm{test}}
813
                  1: \mathbb{U} \leftarrow \{1, \dots, m_t\}, \quad \mathbb{V} \leftarrow \{1, \dots, m_s\}
814
                  2: for i in \mathbb{U} do
                                                                                                                                         \triangleright for each train configuration \omega_i
815
                              best_score \leftarrow -\infty, best_j \leftarrow null
                  3:
816
                  4:
                              \boldsymbol{\omega}_i \leftarrow \boldsymbol{\Omega}_{\text{train}}[\mathcal{A}, i]
817
                              for j in \mathbb{V} do
                  5:
                                                                                                                                          \triangleright find best test configuration \omega_i
818
                  6:
                                     \boldsymbol{\omega}_j \leftarrow \boldsymbol{\Omega}_{\text{test}}[\mathcal{A}, j]
819
                                     s \leftarrow \text{SCORE}\left(\boldsymbol{\omega}_i, \boldsymbol{\omega}_i, \theta\right)
                  7:
                  8:
                                     if s > \text{best\_score then}
820
                  9:
                                            best\_score \leftarrow s, best\_j \leftarrow j
821
                10:
                                     end if
822
                11:
                              end for
823
                12:
                              M \leftarrow \text{PAIR\_MAPPING} (\Omega_{\text{train}}[:, i], \Omega_{\text{test}}[:, \text{best\_j}])
824
                              for p = 1 to N do
                13:
825
                14:
                                     \Omega_{\text{test}}[p, \text{best}_j] \leftarrow M(\Omega_{\text{test}}[p, \text{best}_j])
826
                15:
                              end for
827
                16:
                              Remove i from \mathbb{U}, remove best_j from \mathbb{V}
828
                17: end for
829
                18: return \Omega_{\mathrm{test}}
830
                19: function PAIR_MAPPING(\omega_i, \omega_j)
831
                20:
                              n_i \leftarrow \|\boldsymbol{\omega}_i\|_{\infty}, \quad n_i \leftarrow \|\boldsymbol{\omega}_i\|_{\infty}
832
                              for p = 1 to N do
                                                                                                                               \triangleright build confusion matrix C \in \mathbb{N}^{n_i \times n_j}
                21:
833
                22:
                                     C[\boldsymbol{\omega}_i[p], \boldsymbol{\omega}_j[p]] += 1
834
                23:
                              end for
835
                24:
                              Construct two-walk Laplacian L_{
m tw}
836
                25:
                              \mathcal{F} \leftarrow \text{Fiedler vector of } \boldsymbol{L}_{\text{tw}}
                26:
                              Split \mathcal{F} \to (\mathcal{F}_i \in \mathbb{R}^{n_i}, \, \mathcal{F}_j \in \mathbb{R}^{n_j})
837
                              \pi_i \leftarrow \operatorname{argsort}(\mathcal{F}_i), \ \pi_j \leftarrow \operatorname{argsort}(\mathcal{F}_j)
                27:
838
                28:
                              return mapping k \mapsto \pi_i[\pi_i^{-1}(k)] for k = 1, \dots, \min(n_i, n_i)
839
840
                29: end function
```

MLP (3LP), and three neural tabular architectures (TabNet, TabTransformer, FT-Transformer) in three modes: baseline, static configuration concatenation (GC), and attention-based fusion (GMC).

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

The sole exception is TabTransformer on Boston Housing, where GC yields only a marginal  $R^2$  increase (0.811 $\rightarrow$ 0.813), but GMC degrades it (to 0.671), suggesting that attention-based fusion may disrupt already well-structured features in this case.

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

#### F.2 Additional qualitative evaluation of configurations

In Section 4.3 we provided the embedding of MNIST digits using UMAP and SG-t-SNE (Fig. 8a). Here we provide the missing illustration of embedding with PCA and autoencoder (AE) in Fig. 12. As expected, they do not provide representation with clusters as separated as the former two methods.

With the final figure (Fig. 13) we visualize predicted vs. actual values from the tabular baselines on DSNI, filling in what is missing from Fig. 6.

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

| Dataset  | BHouse       |              | MNIST        |              | CIFAR10      |              |
|----------|--------------|--------------|--------------|--------------|--------------|--------------|
| Metric   | MSE ↓        | $R^2$        | CE↓          | Acc          | CE↓          | Acc          |
| 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  | 0.017        | 0.909        | 0.028        | 0.993        | 0.220        | 0.949        |
| 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 | <u>0.012</u> | <u>0.936</u> | <u>0.017</u> | <u>0.995</u> | <u>0.077</u> | <u>0.978</u> |
| TabT     | 0.035        | 0.811        | 0.192        | 0.980        | 1.028        | 0.706        |
| TabT+GC  | 0.035        | 0.813        | 0.040        | 0.993        | 1.049        | 0.704        |
| TabT+GMC | 0.061        | 0.671        | 0.018        | 0.994        | 0.458        | 0.911        |
| 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  | 0.026        | 0.860        | 0.018        | <u>0.995</u> | 0.157        | 0.955        |

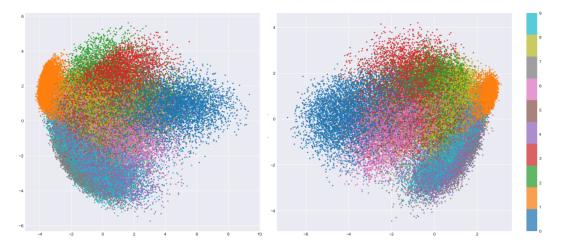


Figure 12: Illustration of 2D embeddings learned by PCA (left) and AE (right) on MNIST.

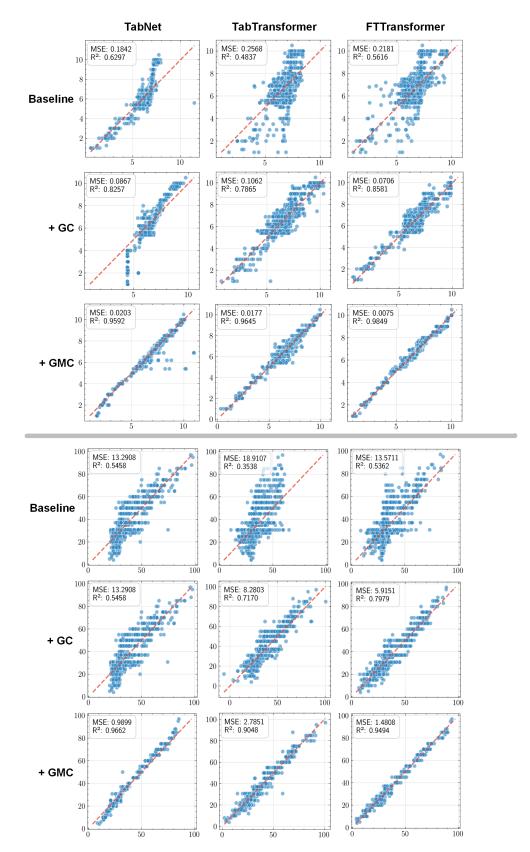


Figure 13: Illustration of the regression performance improvement example in TabNet, TabTransformer and FT-Transformer by adding GC or GMC. Each plots predicted vs. actual value.