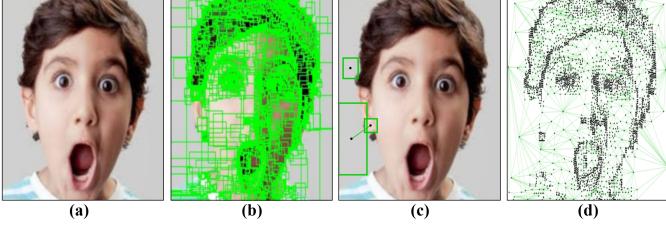


736 We provide more details and results about our work in the
 737 appendices. Here are the contents:

- 738 • Appendix A: Details of Granular Ball Represent.
- 739 • Appendix B: Additional Experiment Results.

740 A Details of Granular Ball Represent

741 In this section, we present in more detail the whole process
 742 of Granular Ball Represent (GBR) for expression images.



743 Figure 7: Process decomposition of GBR . (a) original image, (b) the intermediate process figure of GBR, (c) the simplification of (b),
 744 (d) the final result of GBR.

745 A.1 Fundamental idea of GBR

746 Given an image of expression x^i , by using Algorithm 1
 747 [Xia et al., 2023], we can obtain the GB_list , which is a list
 748 containing several matrices. Previously we have been talking
 749 about the field of granular balls, but the specific implemen-
 750 tation is in the form of granular matrix. This is due to in
 751 two-dimensional surface, the size of images are determined
 752 by x, y, while the granular ball can only rely on the radius
 753 to determine the size of GBR. In [Xia et al., 2023], while in-
 754 novatively applying GBR to images, Xia et al. proposed the
 755 granular moment as a unit of representation which is more
 756 suitable for image data. The GB_list represents the visual-
 757 ization on x^i as shown in Fig. 7b. Each node in the graph x^g
 758 corresponding to x^i is centered by the corresponding matrix
 759 element within GB_list . For each edge of x^g , it consists of
 760 the line joining the centroids of any two matrices with
 761 overlapping regions, as shown in Fig. 7c.

762 A.2 Algorithm of GBR

763 A.3 Features of Node and Edge

764 We tried to design more features related to spatial informa-
 765 tion from the perspective of edge points, expecting to provide
 766 more spatial information to the model.

767 **Features of nodes.** This part of the features consists of the
 768 position coordinates (x, y) of the Node, the variance Var_k ,
 769 mean $Mean_k$, skewness $Skew_k$, maximum Max_k and min-
 770 imum Min_k of the $k - th$ matrix region, and the entropy
 771 Ent_k which can reflect the complexity of the information in
 772 the matrix region, which can be formulated as:

$$773 Fea_n = (x, y, Var_k, Mean_k, Skew_k, Max_k, Min_k, Ent_k), \quad (11)$$

774 where Fea_n denotes the features of node.

775 **Features of edges.** This part of the features consists of the

Algorithm 1 algorithm of GBR

Input: an image of expression: x^i

Parameter: P_{thr} , thr_1 , Var_{thr}

Output: GB_list

```

1: Computing the gradient of each pixel  $p_j$  in  $x^i$ :  $x_{grad}^i$ 
2: Set label  $l(p_j) = 0$  for each pixel  $p_j$  in  $x_{grad}^i$ 
3: Set  $x_{grad}^i = x_{grad}^i.sort$ 
4: repeat
5:   /* Select the center point with minimal gradient.*/
6:   if  $\exists l(p_j) == 0, p_j \in x_{grad}^i$  then
7:      $[grad_{min}, c_k] = x_{grad}.min$ 
8:   end if
9:   /* Rectangular region searching,  $(2 * r_x + 1, 2 * r_y + 1)$  represents (width, length) and  $f(p_j)$  represents the
10:  gray value of  $p_j$ . */
11:  Set  $r_x = 0, r_y = 0$ 
12:  repeat
13:     $(r_x + 1 or r_y + 1) \rightarrow Region R_k$ 
14:     $Purity_k = 1 - \frac{\sum(f(p_j) - f(c_k)) > thr_1}{(2 * r_x + 1)(2 * r_y + 1)}, j \in R_k$ 
15:    Compute the variance of  $R_k$ :  $Var_k$ 
16:    if  $Purity_k < P_{thr} \parallel Var_k > Var_{thr}$  then
17:      if  $r_x(r_y)$  plus one in this iteration. then
18:         $r_x(r_y) = 1(r_y - 1)$  and  $r_x(r_y)$  stop increasing
           in later iterations.
19:      end if
20:    end if
21:    Set  $P_{thr}* = 1.005$ 
22:  until  $r_x \& r_y$  stop increasing.
23:  /* Update the graph.*/
24:   $GB_k = (c_k, r_x, r_y, Purity_k, Var_k, Mean_k)$ 
25:  Add  $GB_k$  to  $GB\_list$ .
26:   $l(p_j) = 1$  for  $p_j$  in  $R_k$ 
27: until  $l(p_j) == 1$  for each pixel  $p_j$  in  $x_{grad}$ 
28: return  $GB\_list$ 

```

length l of the edge, the positional coordinates of the two end
 773 nodes (x_1, y_1, x_2, y_2) , and the area s of the overlap region,
 774 which can be formulated as:
 775

$$776 Fea_e = (l, x_1, y_1, x_2, y_2, s), \quad (12)$$

777 where Fea_e denotes the features of edge.

778 A.4 Details of Graph Down Sample

779 The graph generated by GBR contains a large number of
 780 nodes, some of which are located in the hair, clothing, and
 781 other non-expression related areas, which will import noise
 782 into the model. This will reduce the accuracy of the model
 783 if it is not processed. To address this problem, we perform
 784 graph downsampling in GBR to retain spatial information
 785 while minimizing the number of noise points. As shown in
 786 Fig. 8, the number of noise points is drastically reduced when
 787 10% of nodes is retained. At the same time, the AGM and
 788 L_{BA} modules constructed by our method enable the visual
 789 representation to guide the attention region of the spatial rep-
 790 resentation, which further weakens the negative impact of ir-
 791 relevant regions.

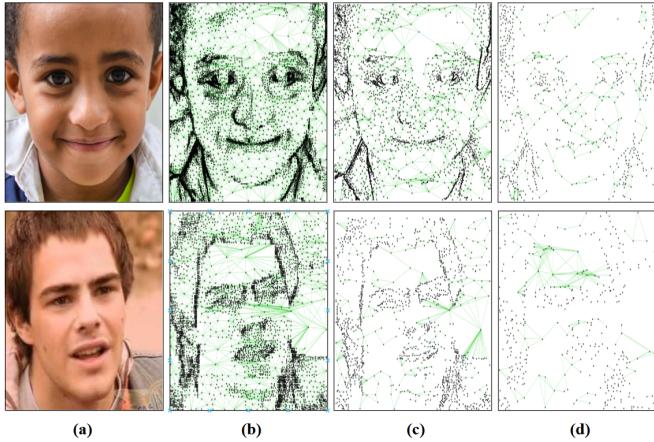


Figure 8: Different downsampling ratios of GDS. (a) original image, (b) Retain 100% of nodes, (c) Retain 50% of nodes, (d) Retain 10% of nodes.

The number of nodes contained in x^g is closely related to the resolution of x^i as well as the tonal complexity. And high resolution images are very common these days. The downsample operation of [Xia *et al.*, 2023] is only dedicated to reduce the nodes with edge connections. In order to effectively reduce the redundant nodes in x^g , we design a downsample operation for edge-connectionless nodes, the basic idea is that for j -th node n_j in x^g which has no edge-connections, the node merging is performed by taking the nearest node n_k , where the distance between the two points needs to be less than a threshold value of d . The two downsample operations are used together in the process of x^g generation. The number of nodes per graph is reduced to within the number N . The value of d is determined by the resolution ($W \times H$) as well as N , which can be formulated as:

$$d = \frac{W}{(\sqrt{N}-1)}, \text{ when } W = H \quad (13)$$

where d is the threshold.

For easier explanation, Eq. 13 simplifies the issue by assuming that W and H are equal. The d derived from Eq. 13 is a limiting case, and it is impossible to have a situation where the distance between all nodes is greater than d and the number of summarized points is still greater than N . d ensures that the total number of nodes after final downsample is within the range N . Our method selects the merge of the nearest nodes, which greatly maintains the quality of the graph data while reducing the number of nodes. The following is the same theoretical reasoning when W and H are not equal, which can be formulated as:

$$\begin{cases} n_w = \frac{n_h \times W}{H} \\ f(n_h) : (\frac{n_h \times W}{H}) \times n_h <= N \\ n_h = \operatorname{argmax}(f(n_h)) \end{cases} \quad (14)$$

$$d = \max(\frac{W}{n_w}, \frac{H}{n_h}) \quad (15)$$

B Additional Experiment Results

B.1 Ablation of Others Setting

In this subsection, we provide ablation studies for the remaining model parameter settings.

downsampling ratios of GDS	Accuracy (%)	
	CAER-S	CK+
100%	92.06%	98.86%
50%	92.76%	99.43%
10%	93.12%	100.00%
5%	92.67%	99.39%

Table 5: The ablation study of different downsampling ratios in the process of GDS. Evaluating on the validation set of the databases.

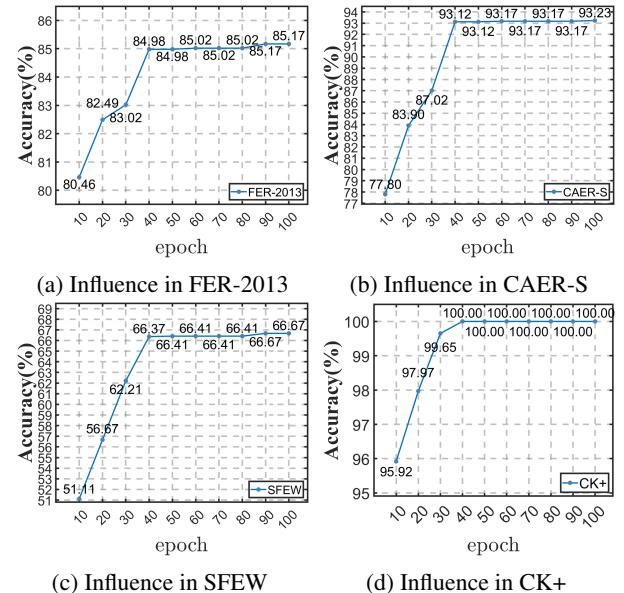


Figure 9: Ablation studies for the different epoch of training on validation set.

Influence of different downsampling ratios. We evaluate the recognition performance of the proposed method with the different downsampling ratios in the process of GDS, as shown in Table 5. The graphs generated by GBR contain spatial information that images do not have. However, the introduced spatial information contains not only valid spatial information but also noise spatial information. We set the appropriate downsampling rate to eliminate the noise as much as possible while introducing the spatial information. Specifically, we can observe that the model achieves the best results with the downsampling rate set to 10%.

Influence of the epoch of training. We evaluate the recognition performance of the proposed method with the different epochs in the process of training, as shown in Fig. 9. Specifically, we can observe that the proposed method achieves a nice recognition accuracy before 40 epochs of training. After that, with the epoch increasing, the recognition accuracy can hardly be greatly improved but the time required for its training increases dramatically.

Influence of the batch size of training. We evaluate the recognition performance of the proposed method with the different batch sizes in the process of training, as shown in

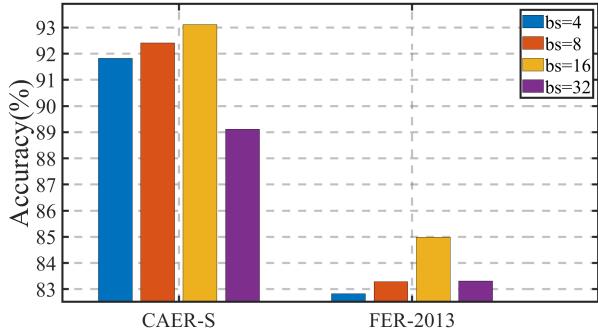


Figure 10: Ablation studies for the different batch size of training on the CAER-S and FER-2013 databases. Note that bs means the batch size of training.

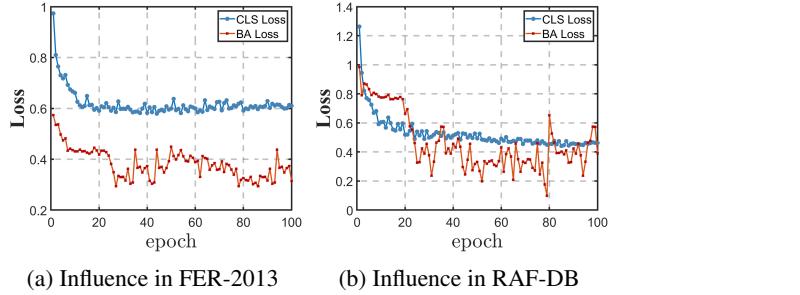


Figure 11: Experiments of convergence analysis on different databases.

845 Fig. 9. We can note that the proposed method achieves the
 846 best accuracy when the batch size is set to 16. Too small a
 847 batch size causes the model to focus too much on the differ-
 848 ences between batches of samples during training to update
 849 the training parameters frequently, which makes the gradient
 850 more stochastic. Increasing the batch size helps with con-
 851 vergence stability, but the generalization performance of the
 852 model decreases as the batch size increases.

853 B.2 Convergence Analysis

854 In this subsection, convergence analysis experiments are
 855 performed on the databases used for our experiments. As
 856 shown in Fig. 11, the increasing number of epochs leads to
 857 a gradual convergence of all loss functions, ultimately reach-
 858 ing a stable state. This clear observation serves as compelling
 859 evidence for the stability and effectiveness of our proposed
 860 model.

861 B.3 More Visualization

862 **2D feature visualization.** We use t-SNE to visualize F_u of
 863 CS-SBF on different test databases on the 2D space, respec-
 864 tively, as shown in Fig. 12. With t-SNE visualization, we can
 865 conclude that the proposed method has good classification re-
 866 sults not only within the network. Visualized on the 2D space,
 867 we can intuitively see that F_u can effectively reduce the intra-
 868 class differences and enhance the inter-class separability of
 869 different expressions.

870 **More Attention visualization.** In this subsection, we pro-

vide more attention visualization of different expression im-
 ages, as shown in Fig. 13.

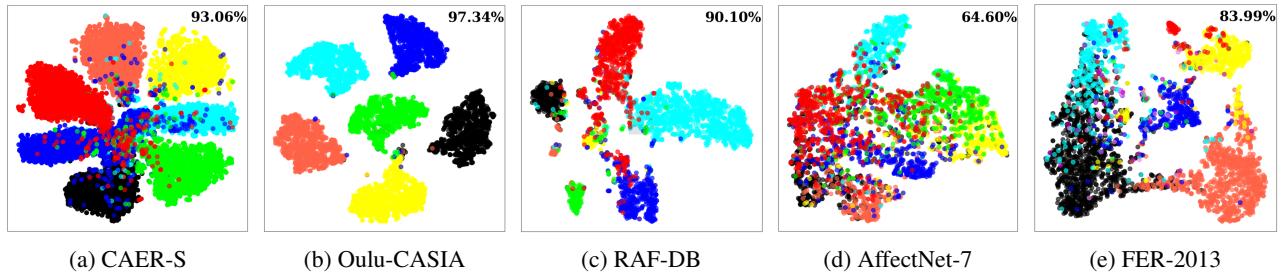


Figure 12: Visualization of fusion features F_u using t-SNE.

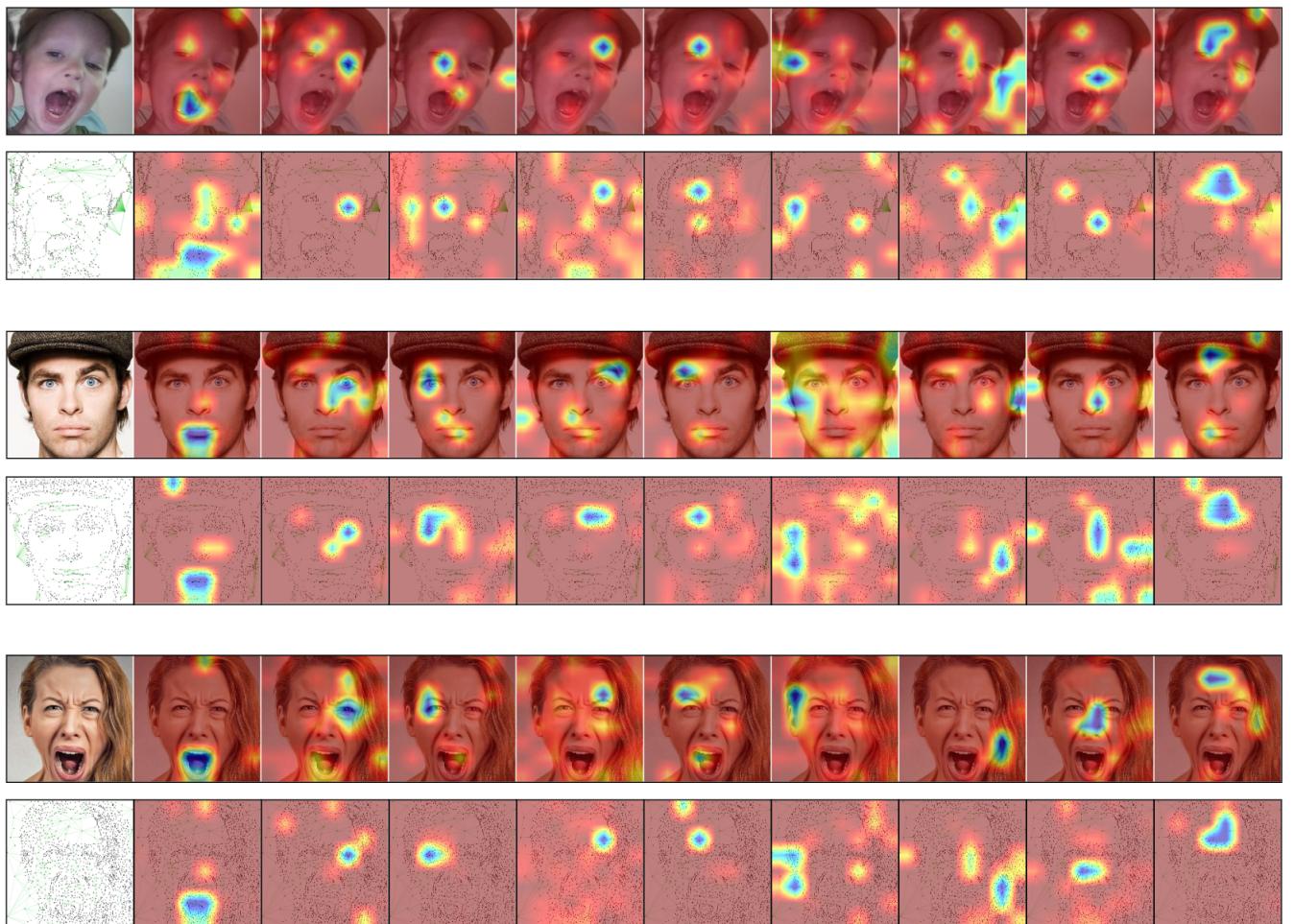


Figure 13: Visualization of component representation features C^w and C^s by CAM on more samples.