

# Clustering-based Feature Representation Learning for Oracle Bone Inscriptions Detection

Ye Tao<sup>1</sup>, Xinran Fu<sup>1</sup>, Honglin Pang<sup>1</sup>, Xi Yang<sup>1,3,4\*</sup>, Chuntao Li<sup>2,4\*</sup>

<sup>1</sup>School of Artificial Intelligence, Jilin University, Qianjin Street,  
Changchun, 130000, Jilin, China.

<sup>2</sup>School of Archaeology, Jilin University, Qianjin Street, Changchun,  
130000, Jilin, China.

<sup>3</sup>Engineering Research Center of Knowledge-Driven Human-Machine  
Intelligence, MoE, Changchun, 130000, Jilin, China.

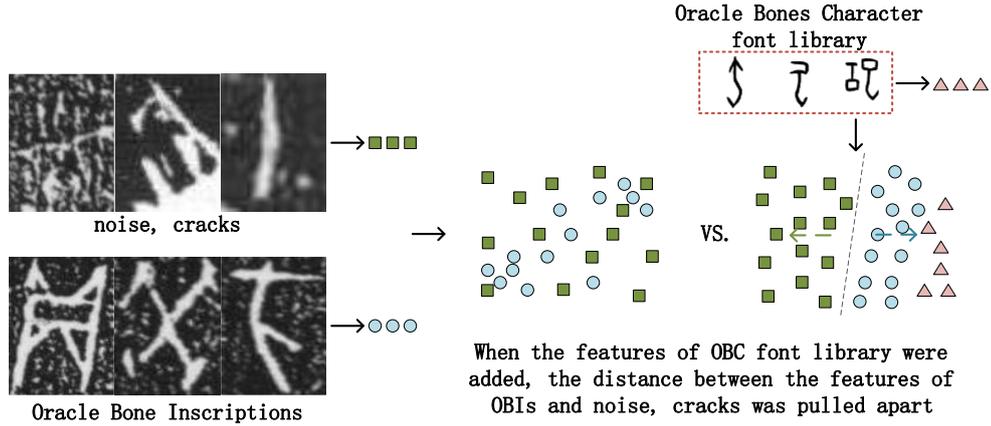
<sup>4</sup>Key Laboratory of Ancient Chinese Script, Cultural Relics and  
Artificial Intelligence, Jilin University, Changchun, 130000, Jilin, China.

\*Corresponding author(s). E-mail(s): [yangxi21@jlu.edu.cn](mailto:yangxi21@jlu.edu.cn);  
[lct33@jlu.edu.cn](mailto:lct33@jlu.edu.cn);

Contributing authors: [taoye23@mails.jlu.edu.cn](mailto:taoye23@mails.jlu.edu.cn);  
[fuXR21@mails.jlu.edu.cn](mailto:fuXR21@mails.jlu.edu.cn); [panghl22@mails.jlu.edu.cn](mailto:panghl22@mails.jlu.edu.cn);

## Abstract

Oracle Bone Inscriptions (OBIs), play a crucial role in understanding ancient Chinese civilization. The automated detection of OBIs from rubbing images represents a fundamental yet challenging task in digital archaeology, primarily due to various degradation factors including noise and cracks that limit the effectiveness of conventional detection networks. To address these challenges, we propose a novel clustering-based feature space representation learning method. Our approach uniquely leverages the Oracle Bones Character (OBC) font library dataset as prior knowledge to enhance feature extraction in the detection network through clustering-based representation learning. The method incorporates a specialized loss function derived from clustering results to optimize feature representation, which is then integrated into the total network loss. We validate the effectiveness of our method by conducting experiments on two OBIs detection dataset using three mainstream detection frameworks: Faster R-CNN, DETR, and Sparse R-CNN. Through extensive experimentation, all frameworks demonstrate significant performance improvements.



**Fig. 1:** Idea of our proposed method. The features of OBIs and noise, cracks were originally mixed together, and when the features of OBC font library were added, the distance between them was pulled apart.

## 1 Introduction

Oracle Bone Inscriptions (OBIs), as the earliest mature writing system in East Asia, represent an invaluable cultural heritage that provides crucial insights into ancient Chinese civilization. Rubbing images, created by placing paper over raised, incised, or textured surfaces and rubbing it with a colored substance, have been used to preserve and study these inscriptions without damaging the original artifacts. Hence the automated detection of OBIs from rubbing images serves as a fundamental prerequisite for character decoding and subsequent historical research, contributing significantly to various fields including traditional Chinese philosophy, astronomy, calendar studies, and historical geography. While traditional methods relied heavily on manual expertise, which was both time-consuming and resource-intensive, modern deep learning approaches have revolutionized this process by offering more efficient and cost-effective solutions. [1–3] However, the detection of OBIs from rubbing images presents several unique challenges: (1) The rubbing images are often contaminated with substantial background noise. (2) The presence of cracks in the rubbing images, which share similar textural properties with OBIs, makes differentiation particularly challenging.

To address these limitations, we propose a novel clustering-based feature space representation learning method that leverages the Oracle Bones Character (OBC) font library dataset [4] as prior knowledge, as shown in Fig. 1. Our approach is founded on the intuition that OBIs and non-OBIs (such as cracks and noise) should occupy distinct regions in the deep feature space. By utilizing the OBC font library as a clean, expert-curated reference point, we guide the model to learn more discriminative feature representations that better distinguish between authentic characters and artifacts.

Previous research in this domain has primarily developed along three key trajectories. The first line of work focuses on OBIs detection.

Previous work has focused on applying deep learning networks to OBIs detection and proposing improved modules for OBIs dataset difficulties. Wang et al. [5] proposed a dynamic augmentation algorithm based on font dataset and an identification auxiliary detection algorithm, and improved the accuracy of compound graphs detection. Xing et al. [6] used the mainstream object detection models to carry out experiments on an OBIs detection dataset and proposed a data augmentation algorithm based on several kinds of noises extracted from the OBIs. Chen et al. [7] proposed a method where shape-adaptive Gaussian kernels are employed to represent the spatial regions of different OBIs. Fu et al. proposed two algorithms to improve the performance of OBIs detection, namely feature fusion based on cross-attention mechanism method [8] and pseudo-category labels prediction method [9]. The above methods can be mainly divided into two categories. The first is to increase the number of difficult samples such as noise and compound graphs utilizing data augmentation to improve the detection effect on difficult samples. The second type is to improve the detection effect by introducing prior knowledge related to OBIs location information, which is obtained through a trained supervised model in advance. The differences between the methods we proposed and others are as follows: (1) The method we proposed is end-to-end, and the introduction of prior knowledge does not need to train other networks in advance. (2) Our method lengthens the distance between OBIs and non-OBIs at the feature space, which improves the detection model’s ability to distinguish them from the data structure, thus improving the detection effect.

A second significant research stream addresses text detection challenges. Previous works in scene text detection can be broadly categorized into regression-based and segmentation-based methods.

Regression-based Methods use the anchor to generate lots of candidate boxes, then NMS[10, 11] is used to obtain the detection results. Firstly, there are methods inspired by object detection. At this stage, the scene text detection method directly locates text by modifying the region proposal of the object detector and the bounding box regression module. TextBoxes [12] define the default box as a quadrilateral with different aspect ratio specifications and use SSD [13] to adapt to different directions and aspect ratios of the text. TextBoxes++ [14] extended horizontal text detection to arbitrary orientations. EAST [15] simplified the detection pipeline by directly predicting text coordinates and angles. Secondly, there are methods based on sub-text components. CTPN [16] introduced vertical anchor regression to localize text lines accurately. SegLink [17] extends CTPN by considering multi-directional links between fragments. There are other methods to improve specific problems. LOMO [18] addressed challenging cases such as long text through iterative refinement, while SBD [19] improved robustness by discretizing quadrilateral boxes into key edges.

Segmentation-based Methods approach text detection through pixel-level semantic segmentation. The Mask-TextSpotter series is An End-to-End Trainable Neural Network for Spotting Text with Arbitrary Shapes. The V1 version[20] used Mask R-CNN as baseline, and can perform end-to-end text spotting. Its mask branch can generate the text instance segmentation maps and the character segmentation maps. The V2 version[21] proposed a spatial attention module to enhance the performance and universality. The V3 version[22] adopts a Segmentation Proposal Network (SPN)

instead of the region proposal network (RPN), which gives accurate representations of arbitrary-shape proposals. To address the difficulty when recognize curved, irregular text, CRAFT [23] detected individual characters and combined them using affinity scores, but it needs complex training. DBNet [24] introduced differentiable binarization to adaptively set thresholds, reducing dependency on post-processing.

The last key research direction centers on cluster-based representation learning. Self-supervised learning [25–30] means that for unlabeled data, the pretext task is designed to mine its own representation features as supervised information to improve the feature extraction ability of the model. In this field, clustering-based representation learning methods have emerged as a particularly promising approach. DEC [31] operates by iteratively optimizing a clustering objective based on KL divergence, using a self-training target distribution. JULE [32] integrates agglomerative clustering with CNNs and frames this combination as a recurrent process. DeepCluster [33] iteratively clusters deep features and utilizes the resulting cluster assignments as pseudo-labels to train the convolutional neural network. DeeperCluster [34] is designed to handle large volumes of uncurated data. Tuo et al. proposed a clustering-based supervised learning scheme for point cloud analysis [35]. They conducted within-class clustering to learn an appropriate point embedding space that is aware of both discriminative semantics and challenging variations. In order to solve the problem that autoencoder is sensitive to noise data in multi-view clustering research, Fatemeh Daneshfar et al. proposed an Elastic Deep Multi-view Autoencoder with Diversity Embedding (EDMVAE-DE) method [36], which has adaptable elastic loss, diversity constraint and Graph regularization to be capable of handling noisy data, it can also make full use of multi-view data.

Specifically, our method integrates both the OBIs detection dataset [37] and the OBC font library dataset during the training process. Through a specialized loss function derived from clustering results, we optimize the feature representation to maximize the separation between true characters and noise while maintaining consistency with the standard forms from the font library. This approach offers several key advantages: 1) It introduces a novel methodology for incorporating expert knowledge through the OBC font library dataset, using pristine character features as anchors for feature space organization. 2) It provides a simple yet effective enhancement that can be readily integrated into existing detection frameworks without requiring complex architectural modifications. 3) It directly addresses the core challenges of OBIs detection by improving the model’s ability to discriminate between authentic characters and various forms of interference.

To validate our approach, we conducted extensive experiments using three mainstream detection frameworks: Faster R-CNN [38], DETR [39], and Sparse R-CNN [40]. The results demonstrate consistent performance improvements across all frameworks, with significant gains in detection accuracy and robustness. Through feature visualization and detailed analysis of detection results, we confirm that our method effectively enhances character feature discrimination while maintaining resilience against various forms of interference. We make several contributions to the field of automated OBIs analysis and digital preservation of ancient writing systems, offering a novel approach

that combines traditional expertise with modern deep learning techniques to achieve superior detection performance.

## 2 Methods

### 2.1 Dataset

There are the introductions of the datasets used in our experiment.

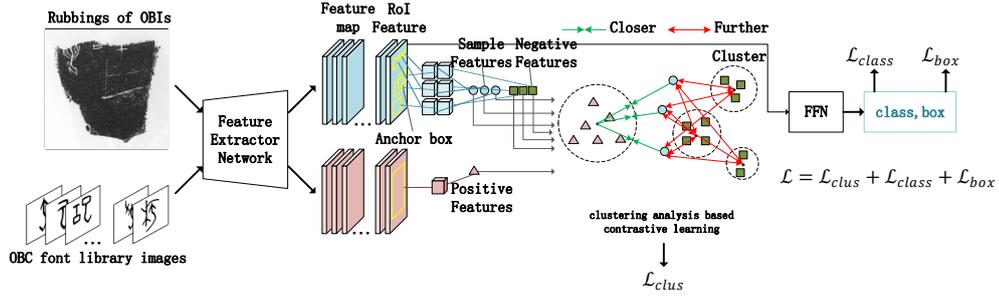
**OBIs Detection Dataset.** We conducted experiments on two datasets. Firstly, we use an open-access OBIs dataset provided by the Key Laboratory of the Ministry of Education for Oracle Information Processing, Anyang Normal University [37]. It has collected 9500 OBI rubbings (up to now) by a high-resolution scanner and then labeled every character with the upper left and lower right coordinates. Notably, this work only focuses on the detection task, so the number of classes in this dataset is all regarded as one.

Secondly, we use OBIMD also provided by the Key Laboratory of the Ministry of Education for Oracle Information Processing, Anyang Normal University [41]. It is the world’s first comprehensive dataset encompassing multi-view information for OBI research, which has collected 10077 OBI rubbings. They are our original input datasets, which are used to train the model for the OBIs detection task. In the following sections, we will refer to the former and the latter as OBIs detection dataset and OBIMD respectively.

**OBC font library Dataset.** Besides, we use another dataset also provided by the Key Laboratory of the Ministry of Education for Oracle Information Processing, Anyang Normal University. It is OBC font library dataset called AYJGW [4], which is used as prior knowledge to enhance feature extraction. With the efforts of authoritative OBI experts, the glyph of every character in AYJGW is confirmed by its inherent meaning. Besides, the font models in the font library are all written by the calligraphy expert. In AYJGW, there are 3,881 character images which have uncontroversial font shapes, and it doesn’t contain any noise. The dataset contains only images but not the annotation information of images, but our work requires the annotation information of its bounding box. Since our task is to detect the OBIs, and each image in this dataset is an independent OBC, the bounding box information of the image can be regarded as the bounding box annotation information of its object. The annotation information of the bounding box for each image in the dataset is formed by the above method.

### 2.2 Overview

The overall pipeline of our method is shown in Fig. 2. Firstly, the OBIs dataset images and the OBC font library images are sent as inputs to the feature extraction network in the detection network at the same time. We will obtain feature maps after the forward propagation of the feature extraction network. We will get anchor boxes based on the feature map respectively, and extract the features of the mapped region of the anchor boxes on the feature map. After the above operations, we will obtain sample features and negative features from the OBIs dataset images, and positive features from the OBC font library images. These features are then flattened into feature vectors, and



**Fig. 2:** The pipeline of our proposed method. Our method takes an OBIs detection dataset and the OBC font library dataset as the inputs. The sample features and negative features are the features extracted in the positive and negative sample boxes on the OBIs images. The positive features are the features extracted in the positive sample boxes on the OBC font library dataset.

we will perform contrastive learning on these feature vectors. In contrastive learning, we will calculate a loss based on the sample and its positive and negative samples, and add this loss to the total loss of the network during training to adjust the model parameters by gradient feedback.

### 2.3 Acquisition of RoI Features

Object detection networks usually first pass the image through a convolutional neural network to extract its feature map. For example, we have applied our method on Faster R-CNN [38], DETR [39], and Sparse R-CNN [40], which all use ResNet[42] as their feature extraction network. First, the OBIs dataset images and the OBC font library images are sent as inputs to the feature extraction network in the detection network at the same time. Let's assume that the image from the OBIs dataset is  $I_1 \in \mathbb{R}^{H_1 \times W_1 \times 3}$ , the image from the OBC font library dataset is  $I_2 \in \mathbb{R}^{H_2 \times W_2 \times 3}$ . We will obtain feature map  $F_1 = E(I_1) \in \mathbb{R}^{H_1 \times W_1 \times C_1}$  and feature map  $F_2 = E(I_2) \in \mathbb{R}^{H_2 \times W_2 \times C_2}$  respectively after the forward propagation of the feature extraction network. We will get anchor boxes based on the feature map, and we will set an Intersection over union(IoU) threshold to judge whether it is a positive anchor box or a negative anchor box according to whether the IoU values of the anchor box and ground truth bounding box are larger or smaller than it. The setting of the threshold is related to whether the original detection model contains sparse box property. If the original detection model contains sparse box property, the threshold will be set relatively small. On the contrary, if the original model does not contain sparse box property, that is, it contains dense box property, then the threshold will be set relatively large. We use the RoI Align layer to extract the features framed by coordinates on the feature map and scale them to a uniform size, and then flatten them to the feature vectors. At this time, the two groups of features to be clustered are the features of OBIs detection dataset images and the OBC font library images, and the next step is to perform cluster analysis.

## 2.4 Clustering-based Contrastive Learning

Our method uses the K-Means clustering algorithm in the training phase to cluster the feature vectors corresponding to the negative anchor boxes in the OBIs detection dataset images and the positive anchor boxes in the OBC font library images respectively. We assume that the number of clusters of the features of the negative anchor boxes in the OBIs dataset is set to be  $N$ , we will get cluster centers  $\{C_N\}$  by using the following formula after K-Means clustering.

$$C_N = \frac{1}{|M_N|} \sum_{x_i \in M_N} x_i \quad (1)$$

where  $M_N$  is the set of feature points belonging to the cluster  $C_N$ .

Similarly, we assume that the number of clusters of the features of the positive anchor boxes in the OBC font library dataset is set to be  $M$ , we will get cluster centers  $\{C_M\}$  after K-Means clustering, then we will take the average at these cluster centers to obtain an average point  $C_M^{mean}$ . We'll conduct a contrastive learning of point and centers of clustering based on the clustering results, in order to achieve intra-cluster compactness and inter-cluster separation. In our method, the samples are the features of the positive anchor boxes in the OBIs images, and their positive samples is  $C_M^{mean}$ , they will form positive sample pairs. Their negative samples are the cluster centers  $\{C_N\}$ . We do this by using the following formula such that each feature point is close to its own cluster center and far away from other cluster centers.

$$\mathcal{L}_{clus}(p_n) = -\log \frac{\exp(\mathbf{p}_n \cdot \mathbf{C}_M^{mean}/\tau)}{\sum_n \exp(\mathbf{p}_n \cdot \mathbf{C}_N/\tau)}, \quad (2)$$

where  $\tau > 0$  is a scalar temperature parameter,  $C_M^{mean}$  refers to the average point of positive sample cluster center of feature point  $p_n$ , which is from cluster centers  $\{C_M\}$ .  $C_N$  refers to the negative sample cluster centers of feature point  $p_n$ , which is from cluster centers  $\{C_N\}$ . The numerator is the similarity of the feature point and its positive sample cluster center, the denominator is the similarity of the feature point and all of its positive sample cluster centers. At training time, we update the parameters in the direction of the descending gradient of the  $\mathcal{L}_{clus}$ , to decrease the distance between  $p_n$  and  $C'$ , while increasing the distance between  $p_n$  and  $C_n$ . The  $\mathcal{L}_{clus}$  will guide the model to train a feature network that makes different features more separate. In the case of our particular task are increasing the distance between the features of OBIs and the features of noise, cracks, and compound graphs. Compared with high-level semantic supervision information, it provides more effective and direct supervision from the aspect of data structure.

## 2.5 Loss function

The input of our method has two image datasets, namely the OBIs detection dataset images and the OBC font library images. During the forward propagation of the network, OBIs images will generate  $\mathcal{L}_{class}$  and  $\mathcal{L}_{box}$ . The total loss function can be

expressed as a weighted sum of  $\mathcal{L}_{\text{class}}$ ,  $\mathcal{L}_{\text{box}}$  and  $\mathcal{L}_{\text{clus}}$ :

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{clus}} + \lambda_2 \mathcal{L}_{\text{class}} + \lambda_3 \mathcal{L}_{\text{box}} \quad (3)$$

The OBC font library images are treated as prior knowledge to be introduced to guide the learning process of the model, which should only be used when calculating the  $\mathcal{L}_{\text{clus}}$ , and other losses they generate are not included in the final total loss. The weight parameters  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$  are set differently according to the characteristics of each model.

## 3 Results

### 3.1 Implementation Details and Evaluation Metrics

We applied our method to three detection models: Faster R-CNN, DETR, and Sparse R-CNN, and they are all implemented based on Python and PyTorch. We trained these models on four NVIDIA RTX A6000. We will introduce the details of experimental implementation for each of the three models.

Faster R-CNN and Sparse R-CNN both used ResNet-50 as feature extraction network, then used Feature Pyramid Networks (FPN)[43] to achieve better feature map fusion. We initialized the parameters of ResNet-50 with weights pre-trained on the ImageNet dataset Torchvision officially provided, and the other parameters were initialized randomly. We employed the SGD optimizer for training Faster R-CNN[44], with a momentum parameter set to 0.9 and a weight decay parameter set to  $1 \times 10^{-4}$ .

DETR used ResNet-101 as the feature extraction network. We initialized the parameters of our DETR model with weights pre-trained on the COCO 2017 dataset Facebook officially provided. DETR and Sparse R-CNN both employed the AdamW[45] optimizer for training, with a weight decay parameter set to  $1 \times 10^{-4}$ .

The remaining parameter settings of the three models are shown in the table 1.

**Table 1:** The setting of learning rate,  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$ , and Batchsize of Faster R-CNN+ours, DETR+ours and Sparse R-CNN+ours.

Model	Learning Rate	$\lambda_1, \lambda_2, \lambda_3$	Batchsize
Ours (Faster R-CNN)	$1 \times 10^{-2}$	1, 1, 1	2
Ours (DETR)	$1 \times 10^{-4}$	1, 5, 1	2
Ours (Sparse R-CNN)	$2.5 \times 10^{-5}$	0.1, 1, 1	4

Before introducing the evaluation metrics of the method, we first introduce the concept of confusion matrix. The sample is a very important concept in the evaluation of object detection methods. Positive sample is the object that needs to be detected, and negative sample is the object that does not need to be detected. When the model is predicting, there are usually four situations: TP, FP, TN, and FN. In object detection, the final bounding boxes predicted by the model are regarded as positive samples, while whether the prediction belongs to TP or FP depends on the value of IoU. IoU

calculates the ratio between the intersection area of the predicted bounding box and the ground truth bounding box and the merging area, and we will set a threshold value, which is 0.5 in this study. If the value of IoU of the predicted bounding box and ground truth bounding box is greater than this threshold, the predicted case of this bounding box is TP. Conversely, if the value of IoU of the predicted bounding box and any ground truth bounding box is less than this threshold, then the prediction of this bounding box is FP. If the IoU value of a ground truth bounding box and any predicted bounding box is less than this threshold, it is FN.

We use three evaluation metrics [46] commonly used in object detection models based on confusion matrix, namely Precision, defined as is the proportion of TP in all prediction numbers; Recall, defined as is the proportion of TP in all object numbers; F1-score, defined as is the harmonic mean of precision and recall, which is based on the comprehensive evaluation model of precision and recall.

### 3.2 Main results

To prove the performance improvement of our method for detection models, we conducted a series of comparative experiments on the two OBIs detection datasets mentioned in Section 2.1, Precision, Recall, and F1-score were used to evaluate and compare the performance of the Faster R-CNN, DETR, DENO[47] and Sparse R-CNN model and the improved Faster R-CNN, DETR, DENO, and Sparse R-CNN model, as shown in Table 2 and Table 3. The model without special marking means that its backbone is ResNet.

**Table 2:** Comparison results on OBIs detection dataset [37] between existing models and our proposed method. AP is averaged over multiple IoU values. Specifically we use 10 IoU thresholds of .50:.05:.95. AP<sub>50</sub> is computed the average precision at a single IoU of 0.50. AR<sub>50</sub> is computed the average recall at a single IoU of 0.50.

Model	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>50</sub>	F1-Score <sub>50</sub>
Faster R-CNN-R50	47.0	83.7	47.9	89.1	86.3
Faster R-CNN-Swin	48.0	86.8	48.8	61.1	71.7
DETR-R101	35.7	75.8	28.4	88.2	81.5
DENO-Swin	39.9	78.4	36.1	91.7	84.5
Sparse R-CNN-R50	36.5	72.3	33.1	89.0	79.7
Sparse R-CNN-Swin	36.5	71.4	33.3	45.9	55.9
<b>Ours</b> (Faster R-CNN-R50)	48.1(+1.1)	86.4(+2.7)	48.7(+0.8)	92.5(+3.4)	89.3(+3.0)
<b>Ours</b> (Faster R-CNN-Swin)	47.8(-0.2)	86.6(-0.2)	48.2(-0.6)	62.4(+1.3)	72.5(+0.8)
<b>Ours</b> (DETR-R101)	36.7(+1.0)	77.0(+1.2)	30.7(+2.3)	88.3(+0.1)	82.3(+0.8)
<b>Ours</b> (DENO-Swin)	41.7(+1.8)	81.3(+2.9)	38.6(+2.5)	92.5(+0.8)	86.5(+2.0)
<b>Ours</b> (Sparse R-CNN-R50)	37.2(+0.7)	73.4(+1.1)	34.1(+1.0)	89.0	80.4(+0.7)
<b>Ours</b> (Sparse R-CNN-Swin)	37.2(+0.7)	72.4(+1.0)	34.6(+1.3)	47.0(+1.1)	57.0(+1.1)

**Table 3:** Comparison results on OBIMD [41] between existing models and our proposed method.

Model	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>50</sub>	F1-Score <sub>50</sub>
Faster R-CNN-R50	33.2	70.9	26.2	79.8	75.1
DETR-R101	29.3	66.4	21.0	75.7	70.7
Sparse R-CNN-R50	29.2	60.1	23.3	77.7	67.7
<b>Ours</b> (Faster R-CNN-R50)	34.3(+1.1)	73.4(+2.5)	26.6(+0.4)	81.4(+1.6)	77.2(+2.1)
<b>Ours</b> (DETR-R101)	30.0(+0.7)	68.0(+1.6)	22.0(+1.0)	76.1(+0.4)	71.8(+1.1)
<b>Ours</b> (Sparse R-CNN-R50)	30.1(+0.9)	62.6(+2.5)	24.7(+1.4)	78.7(+1.0)	69.7(+2.0)

**Table 4:** Comparison Training time and Memory usage on OBIs detection dataset [37] and OBIMD [41] between existing models and our proposed method.

Dataset	Model	Training time	Memory usage(GiB)
OBIs detection dataset	Faster R-CNN	10h56min	18.77
	<b>Ours</b> (Faster R-CNN)	13h23min	20.57
	DETR	13h12min	10.50
	<b>Ours</b> (DETR)	17h33min	13.01
	Sparse R-CNN	7h32min	12.31
	<b>Ours</b> (Sparse R-CNN)	9h42min	16.14
OBIMD	Faster R-CNN	14h1min	18.89
	<b>Ours</b> (Faster R-CNN)	16h27min	20.90
	DETR	15h8min	10.62
	<b>Ours</b> (DETR)	20h25min	13.23
	Sparse R-CNN	9h35min	12.50
	<b>Ours</b> (Sparse R-CNN)	11h45min	16.24

As shown in Table 2 and Table 3, Considering transformer based approaches like Swin Transformer can provide deeper insights into the advantages of clustering based representation learning, we added the experiment of replacing the backbone of the three models with Swin Transformer [48]. Among them, because DETR replaced the backbone with the Swin Transformer and there was no global pre-training weight, the model failed to converge during training. So we replaced DETR with DINO: DETR with Improved DeNoising Anchor Boxes for End-to-End Object Detection [47]. It is an improved version of DETR and can converge better during training. From the numerical point of view, after the application of our proposed method in these three models, the precision, recall, and F1-score have been improved to some extent compared with the baseline model, only the accuracy rate on Faster R-CNN-Swin has slightly decreased. so our proposed method is meaningful for the improvement of the detection effect of the models.

As shown in Table 4, from the numerical point of view, after the application of our proposed method in these three models, the consumption in terms of both training time and memory usage have increased to a certain extent, but it is all within an acceptable range.

As is well known, there are some high-performance and most popular detection frameworks such as YOLO [49]. However, our method is not suitable for application to models of the YOLO class. Because YOLO ultimately uses global features to predict categories and regression box coordinates, while the method we use takes out each feature in the box to predict categories and regression box coordinates, and our clustering algorithm improves for each feature, it has no effect on global features.

In addition to the comparison of numerical results, we also did visualization experiments, namely the visualization of features and the visualization of detection results.

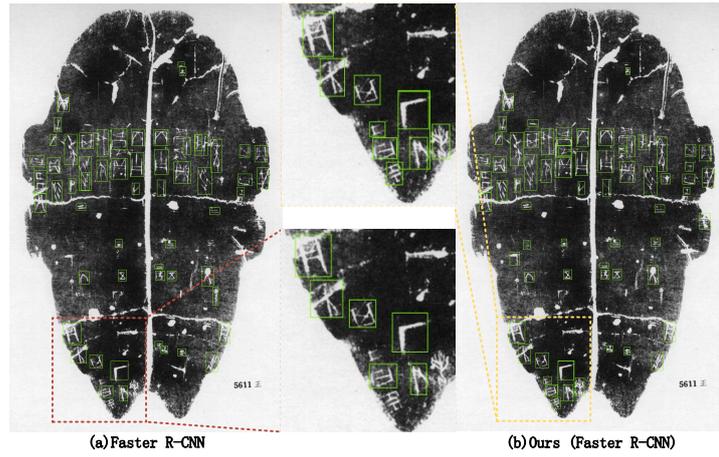
**Detection results.** We visualized the detection results of the three models. As shown in Fig. 3 and Fig. 4, it can be obviously seen that our method has significantly improved the detection effect of the three models.

**Features.** The feature vectors for comparison learning were sample features, positive features, and negative features, which were reduced to two dimensions using t-SNE [50] dimensionality reduction method, and these points were presented in an image. We can use these visualizations to further visualize the role of contrast learning in our method.

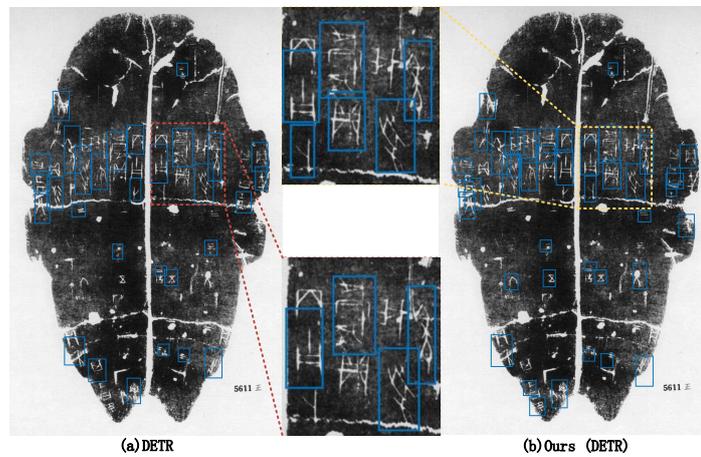
As shown in Fig. 5, we observe the visualization of the evolution process of the Faster R-CNN model using our proposed method first. With the epoch number increased, the sample features, which were initially scattered and unordered, mixed with negative features, are gradually guided by positive features to concentrate at the bottom and linearly separate from negative features. Throughout the process, the number of sample features increases gradually, while the number of negative features decreases. The visualization results of this evolution indicate that the features from OBC font library dataset are effective in guiding the features from OBIs dataset to be away from those of non-characters. Furthermore, as shown in Fig. 6, by directly comparing the feature point visualization image of the baseline Faster R-CNN model and the Faster R-CNN model using our proposed method, it can be seen that the sample features and negative features in the baseline Faster R-CNN model are mixed together and linearly inseparable. In contrast, the sample features and negative features in the Faster R-CNN model using our method are separated and linearly separable. Moreover, the number of sample features in the latter is much greater than that in the former, and the number of negative features in the latter is much smaller than that in the former, which indicates that applying our method can guide the model to output more anchor boxes containing characters, while reducing the output of anchor boxes containing cracks, noise, etc. This has a certain significance for improving the detection performance of the model.

### 3.3 Hyperparameter and Result analysis

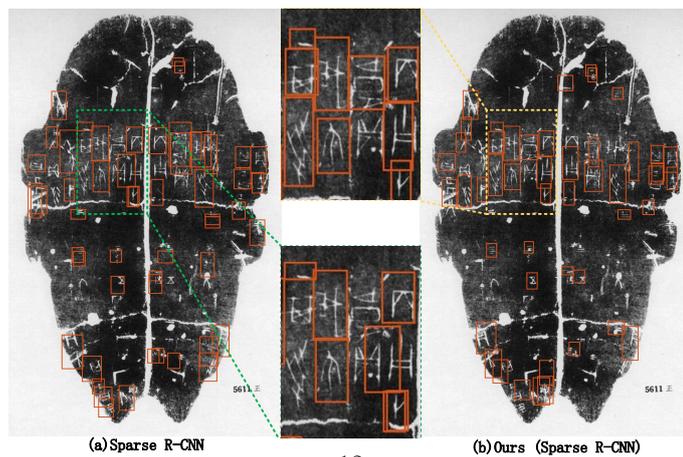
Firstly, we analyze the different effects of the setting of four hyper-parameters, namely the number of OBC from OBC font library dataset used as prior knowledge, the setting of temperature parameter  $\tau$ , the setting of  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$ , and the setting of clustering technique on the model detection effect. Secondly, we analyze whether the results are reproducible.



(a) Comparative detection results on Faster R-CNN

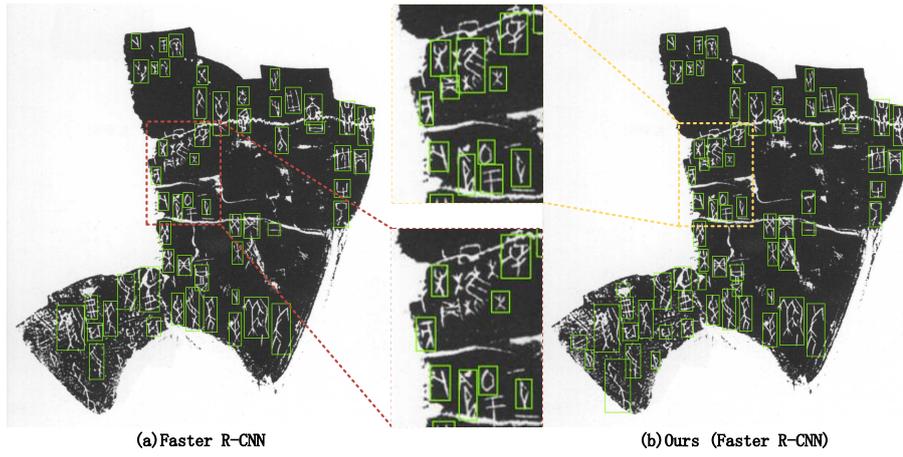


(b) Comparative detection results on DETR

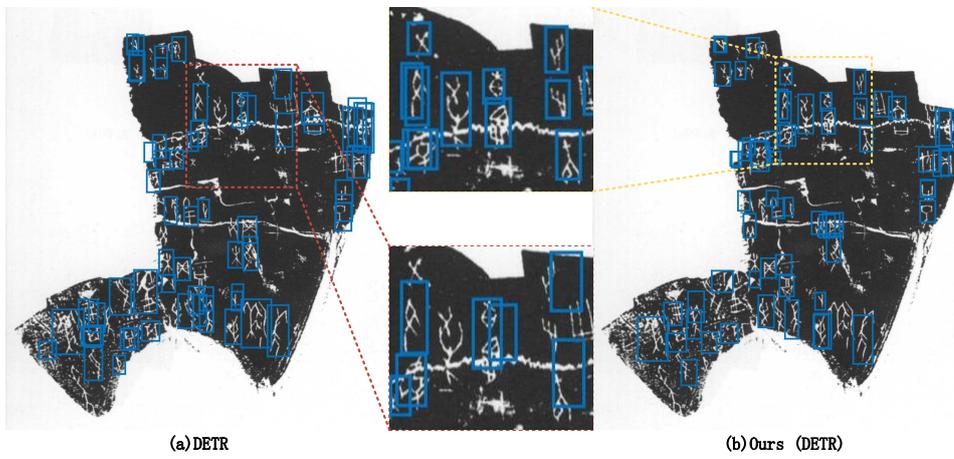


(c) Comparative detection results on Sparse R-CNN

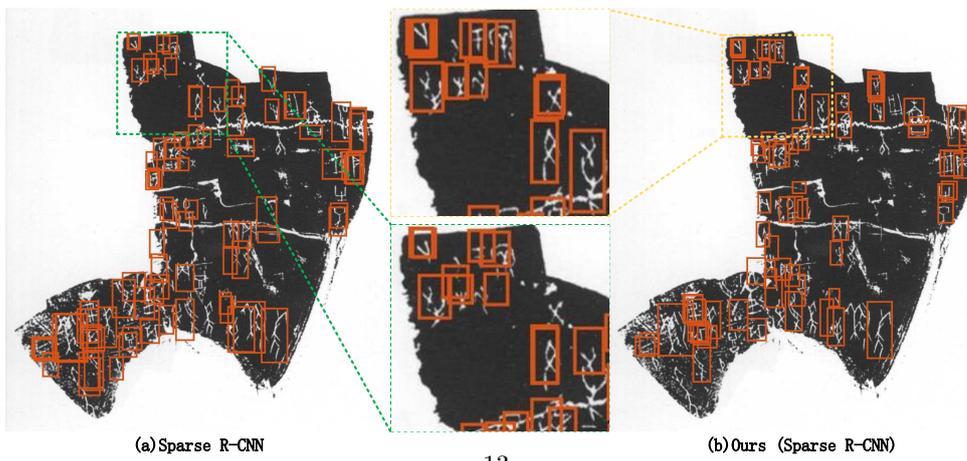
**Fig. 3:** Comparative detection results on three models, the partially enlarged image at the bottom is the baseline model, and the one at the top is with our method.



(a) Comparative detection results on Faster R-CNN

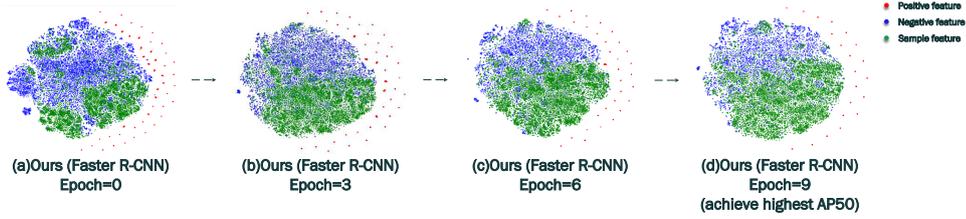


(b) Comparative detection results on DETR

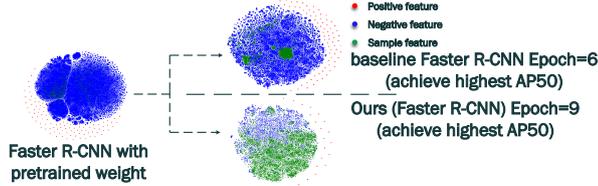


(c) Comparative detection results on Sparse R-CNN

**Fig. 4:** Comparative detection results on three models, the partially enlarged image at the bottom is the baseline model, and the one at the top is with our method.



**Fig. 5:** Comparison of visualization results of the distribution of positive features, negative features and sample features generated from the Faster R-CNN+ours model which uses weight from epoch=0, 3, 6, 9.



**Fig. 6:** Comparison of visualization results of the distribution of positive features, negative features and sample features generated from the baseline Faster R-CNN model which uses weight that reaches the highest  $AP_{50}$  and positive features, negative features and sample features generated from the Faster R-CNN+ours model which uses weight that reaches the highest  $AP_{50}$ .

**Number of OBC used as prior knowledge.** We applied our method on the Faster R-CNN, DETR and Sparse R-CNN models with the number of OBC used as prior knowledge sets of 10, 20, and 50 respectively, and the results obtained are shown in Table 5.

Firstly, we present the experimental results of different numbers of OBC used as prior knowledge on the Faster R-CNN model. When the number of OBC used as prior knowledge is set to 10 and 20, the performance on AP,  $AP_{50}$ , and  $AP_{75}$  is improved compared to the baseline Faster R-CNN model. The improvement effect of the number of OBC used as prior knowledge set to 20 is better than that of the number of OBC used as prior knowledge set to 10. When the number of OBC used as prior knowledge is set to 50, the performance on AP and  $AP_{50}$  is slightly decreased compared to the baseline Faster R-CNN model, and the performance on  $AP_{75}$  is basically the same. Therefore, after the number of OBC used as prior knowledge is set to 50, our method has a negative effect on the Faster R-CNN model, resulting in lower accuracy than the baseline Faster R-CNN model.

We perform a visual analysis of this experimental result by visualizing the distribution of sample features, positive features, and negative features when the number of OBC used as prior knowledge is set to 10, 20, and 50. As shown in Fig. 7, when the number of OBC used as prior knowledge is set to 10 and 20, the positive features

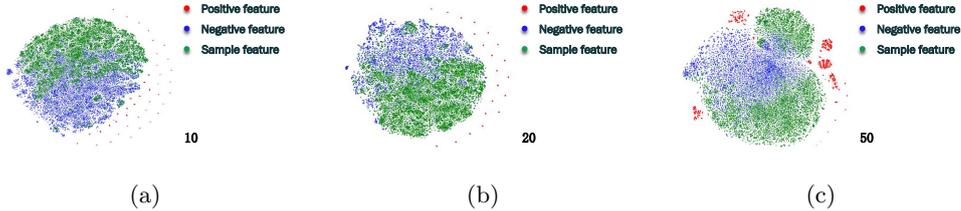
are concentrated on the right side, providing effective guidance for the distribution of sample features, and finally, sample features and negative features can be linearly separated. When the number of OBC used as prior knowledge is set to 50, the positive features are distributed relatively randomly, providing negative guidance for the distribution of sample features, and finally, sample features and negative features cannot be linearly separated.

Secondly, we look at the experimental results of different number of OBC used as prior knowledge on the DETR model. When the number of OBC used as prior knowledge is set to 10, there is a decrease in the performance compared with baseline DETR model on AP, AP<sub>50</sub> and AP<sub>75</sub>. When the number of OBC used as prior knowledge is set to 20 and 50, the effect on AP, AP<sub>50</sub> and AP<sub>75</sub> is improved compared with baseline DETR model. In terms of overall improvement effect, the number of OBC used as prior knowledge set to 50 is better than the set to 20.

The experimental results of different numbers of OBC used as prior knowledge on the Sparse R-CNN model are similar to that for the DETR model. The best setting is 20.

**Table 5:** Effect of number of OBC used as prior knowledge on different models. The performance of the number of OBC used as prior knowledge set to 0 means the baseline performance.

OBC	Faster R-CNN			DETR			Sparse R-CNN		
	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>
0	47.0	83.7	47.9	35.7	75.8	28.4	36.5	72.3	33.1
10	47.5(+0.5)	85.2(+1.5)	49.1(+1.2)	35.2(-0.5)	75.5(-0.3)	27.8(-0.6)	36.3(-0.2)	71.7(-0.6)	33.1
20	48.1(+1.1)	86.4(+2.7)	48.7(+0.8)	36.0(+0.3)	76.7(+0.9)	28.9(+0.5)	37.2(+0.7)	73.4(+1.1)	34.1(+1.0)
50	46.7(-0.3)	83.2(-0.5)	48.0(+0.1)	36.8(+1.1)	76.5(+0.7)	31.1(+2.7)	36.7(+0.2)	72.4(+0.1)	33.4(+0.3)



**Fig. 7:** Comparison of visualization results of the distribution of positive features, negative features and sample features generated from the Faster R-CNN+ours model with the number of OBC used as prior knowledge is set to 10, 20, 50.

**Setting of temperature parameter  $\tau$ .** Our method was experimented on three models and proved to improve the detection performance of all three models compared to the original baseline model. However, it can be observed that the improvement in the sparse detection method is not as significant as that in the dense detection

method. The reason for this result is that in the contrastive learning process of the dense detection method in our method, each sample has more negative samples available for learning, while the number of negative samples corresponding to each sample in the sparse detection method is relatively smaller. In contrastive learning, the temperature plays a role in controlling the strength of penalties on the hard negative samples. Specifically, contrastive loss with small temperature tends to penalize much more on the hardest negative samples. On the other hand, contrastive loss with large temperatures is less sensitive to the hard negative samples. In our method, there is a problem that the number of negative samples corresponding to samples is small in the sparse detection model, so the distinction between positive and negative samples should be enhanced, and the smaller the corresponding temperature should be set, the better the experimental results will be. We did the following experiment.

**Table 6:** Effect of Setting of temperature parameter  $\tau$  on different models. The performance of temperature parameter  $\tau$  set to 0.1 is viewed as the baseline performance.

OBC	$\tau$	Ours (DETR)			Ours (Sparse R-CNN)		
		AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>
10	0.1	35.2	75.5	27.8	36.3	71.7	33.1
10	0.05	36.7(+1.5)	77.0(+1.5)	30.7(+2.9)	37.3(+1.0)	72.8(+1.1)	34.7(+1.6)
10	0.01	36.5(+1.3)	76.1(+0.6)	30.8(+3.0)	36.5(+0.2)	72.0(+0.3)	33.4(+0.3)
10	0.005	34.6(-0.6)	75.1(-0.4)	27.7(-0.1)	36.6(+0.3)	72.2(+0.5)	33.8(+0.7)
20	0.1	36.0	76.7	28.9	37.2	73.4	34.1
20	0.05	37.4(+1.4)	76.8(+0.2)	32.4(+3.5)	36.9(-0.3)	72.7(-0.7)	33.1(-1.0)
20	0.01	36.0	74.9(-1.8)	30.3(+1.4)	37.0(-0.2)	72.9(-0.5)	33.4(-0.7)
20	0.005	35.2(-0.8)	75.2(-1.5)	28.9	36.6(-0.6)	72.1(-1.3)	33.4(-0.7)
50	0.1	36.8	76.5	31.1	36.7	72.4	33.4
50	0.05	35.3(-1.5)	75.8(-0.7)	27.5(-3.6)	35.9(-0.8)	71.4(-1.0)	32.5(-0.9)
50	0.01	38.5(+1.7)	76.8(+0.3)	33.7(+2.6)	37.3(+0.6)	73.0(+0.6)	34.5(+1.1)
50	0.005	34.3(-2.5)	75.3(-1.2)	26.2(-4.9)	36.3(-0.4)	71.7(-0.7)	33.2(-0.2)

From Table 6, it can be seen that when the number of OBC used as prior knowledge is set to 10, the temperature parameter  $\tau$  decreased from 0.1 to 0.05 and 0.01, the performance on AP, AP<sub>50</sub>, and AP<sub>75</sub> is improved compared to the baseline DETR and Sparse R-CNN model. The best setting of the temperature parameter  $\tau$  for the DETR and Sparse R-CNN model are both 0.05 when the number of OBC used as prior knowledge is set to 10. When the number of OBC used as prior knowledge is set to 20, the temperature parameter  $\tau$  decreased from 0.1 to 0.05, the performance on AP, AP<sub>50</sub>, and AP<sub>75</sub> is improved compared to the baseline DETR model. The best setting of the temperature parameter  $\tau$  for the DETR and Sparse R-CNN model are respectively 0.05 and 0.1 when the number of OBC used as prior knowledge is set to 20. When the number of OBC used as prior knowledge is set to 50, the temperature parameter  $\tau$  decreased from 0.1 to 0.01, the performance on AP, AP<sub>50</sub>, and AP<sub>75</sub> is improved compared to the baseline DETR and Sparse R-CNN models. The best setting of the temperature parameter  $\tau$  for the DETR and Sparse R-CNN model are both 0.01 when the number of OBC used as prior knowledge is set to 50. Therefore,

when our method is applied to the sparse detection model, the temperature should be set smaller to compensate for the problem that when the sparse detection models use our method, they will have fewer negative samples corresponding to the samples in the contrastive learning process. This will enable our method to achieve better performance when applied to the sparse detection model.

To further optimize our method performance, we use Bayesian Optimization technique to tune our hyperparameters, and the results obtained were shown in Table 7.

**Table 7:** Use Bayesian Optimization technique to tune our hyperparameters. The M and N are the number of clusters of negative samples and the number of clusters of positive samples mentioned in Section 2.4.

Model	OBC	$\tau$	M,N	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>50</sub>	F1-Score <sub>50</sub>
Faster R-CNN				47.0	83.7	47.9	89.1	86.3
DETR				35.7	75.8	28.4	88.2	81.5
Sparse R-CNN				36.5	72.3	33.1	89.0	79.7
<b>Ours</b> (Faster R-CNN)	27	0.099	63,3	48.5(+1.5)	86.6(+2.9)	48.7(+0.8)	92.6(+1.5)	89.5(+2.2)
<b>Ours</b> (DETR)	40	0.069	41,2	36.7(+1.0)	77.1(+1.3)	31.0(+2.6)	88.5(+0.3)	82.4(+0.9)
<b>Ours</b> (Sparse R-CNN)	34	0.073	48,4	37.3(+0.8)	73.4(+1.1)	34.2(+1.1)	89.0	80.4(+0.7)

**Setting of  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$ .** The total loss function of our method can be expressed as a weighted sum of  $\mathcal{L}_{class}$ ,  $\mathcal{L}_{box}$  and  $\mathcal{L}_{clus}$ :

$$\mathcal{L} = \lambda_1 \mathcal{L}_{clus} + \lambda_2 \mathcal{L}_{class} + \lambda_3 \mathcal{L}_{box} \quad (4)$$

Among them,  $\lambda_2$  and  $\lambda_3$  are set based on the default parameters of the original model, while  $\lambda_1$  is set according to the rule that the loss we add is of the same order of magnitude as the other losses. We applied our method on the Faster R-CNN, DETR and Sparse R-CNN models with different setting of  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$ , and the results obtained were shown in Table 8.

**Table 8:** Effect of different setting of  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$  on different models. The performance of None means the baseline performance.

$\lambda_1, \lambda_2, \lambda_3$	Ours (Faster R-CNN)			$\lambda_1, \lambda_2, \lambda_3$	Ours (DETR)			$\lambda_1, \lambda_2, \lambda_3$	Ours (Sparse R-CNN)		
	AP	AP <sub>50</sub>	AP <sub>75</sub>		AP	AP <sub>50</sub>	AP <sub>75</sub>		AP	AP <sub>50</sub>	AP <sub>75</sub>
1,1,1	47.0	83.7	47.9	1,5,1	35.7	75.8	28.4	0.1,1,1	36.5	72.3	33.1
0.33,0.33,0.33	48.1(+1.1)	86.4(+2.7)	48.7(+0.8)	0.14,0.71,0.14	36.7(+1.0)	77.0(+1.2)	30.7(+2.3)	0.05,0.48,0.48	37.2(+0.7)	72.4(+0.1)	34.6(+1.5)
2,1,1	47.2(+0.2)	85.4(+1.7)	48.3(+0.4)	2,5,1	34.8(-0.9)	75.2(-0.6)	28.3(-0.1)	0.2,1,1	36.8(+0.3)	72.1(-0.3)	34.2(+0.9)

As shown in Table 8, We attempted to set the sum of  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$  to equal 1, or keep the values of  $\lambda_2$ ,  $\lambda_3$  and only increase  $\lambda_1$ . However, the results were not better than keeping  $\lambda_2$  and  $\lambda_3$  in the original model set and then ensuring that  $\lambda_1$  is of the same order of magnitude as them.

**Setting of clustering technique.** We applied our method on the Faster R-CNN, DETR and Sparse R-CNN models with different setting of clustering technique, namely K-Means and DBSCAN, and the results obtained were shown in Table 9.

**Table 9:** Effect of setting of clustering technique on different models. We use the parameters eps=12, min samples=5 and eps=5, min samples=2 and eps=5, min samples=2 for DBSCAN respectively. The performance of the None means the baseline performance.

	Faster R-CNN			DETR			Sparse R-CNN		
	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>
None	47.0	83.7	47.9	35.7	75.8	28.4	36.5	72.3	33.1
K-Means	48.1(+1.1)	86.4(+2.7)	48.7(+0.8)	36.0(+0.3)	76.7(+0.9)	28.9(+0.5)	37.2(+0.7)	73.4(+1.1)	34.1(+1.0)
DBSCAN	46.5(-0.5)	83.3(-0.4)	47.8(-0.1)	30.5(-5.2)	70.4(-5.4)	21.7(-6.7)	36.9(+0.4)	72.2(-0.1)	34.1(+1.0)

As shown in Table 9, after changing the clustering method to DBSCAN, the detection effect of the model has decreased. The reason might be that: DBSCAN marks the points in the low-density area as noise, and our method precisely requires these points determined as noise as negative samples in contrastive learning.

**The reproducibility of result.** In our method, we use random initialization for K-means. Considering the potential sensitivity of K-means to initialization, we conducted three experiments respectively on three models with the same parameters, and the results are as shown in Table 10.

**Table 10:** Effect of different random initializations on different models. We repeated 5 times with different random initializations on different models.

Model	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>50</sub>	F1-Score <sub>50</sub>
Ours (Faster R-CNN)	48.1 ± 0.1	86.0 ± 0.2	49.6 ± 1.1	91.8 ± 0.5	88.8 ± 0.3
Ours (DETR)	36.6 ± 0.2	76.7 ± 0.2	30.5 ± 0.2	88.4 ± 0.3	82.2 ± 0.2
Ours (Sparse R-CNN)	37.2 ± 0.1	73.2 ± 0.1	34.4 ± 0.2	88.9 ± 0.1	80.3 ± 0.1

As shown in Table 10, each experiment was repeated 5 times with different random initializations. We report the mean and standard deviation of the precision and recall. The results are reproducible when rerunning the code.

## 4 Discussion

We propose a novel clustering-based feature space representation learning method to address the unique challenges of the detection of OBIs. The idea of our method is to guide the model to learn more discriminative feature representations by utilizing the OBC font library as a clean, expert-curated reference point. We conducted experiments on three main-stream detection frameworks: Faster R-CNN, DETR, and Sparse R-CNN on two OBIs detection datasets. Both numerical and visual experimental results prove that our method can improve the detection effect of the three models to some extent. Our method offers an idea: how to use font libraries as the prior knowledge of the network to guide the learning of the network and improve the effect of the network, which is generalizable, particularly with the advancements of character databases and font libraries.

**Data availability.** The data are available at <https://jgw.aynu.edu.cn> and <https://www.jgwlbq.org.cn/home> respectively.

**Code availability.** The code are available at <https://github.com/biscuit030/Clustering-based-Feature-Representation-Learning-for-Oracle-Bone-Inscriptions-Detection>

**Acknowledgements.** This work is supported by the Young Scientists Fund of the National Natural Science Foundation of China (Grant No.62206106), and is a phased achievement of the National Social Science Foundation project "Research on Chinese Pre Qin Language and Culture Based on Jinwen Data" (No. 23VRC033), and has also received funding from the interdisciplinary cultivation project for young teachers and students at Jilin University, "Research on Jinwen Data Based on Artificial Intelligence" (No. 2024-JC XK-04), and is funded by the "Ancient Chinese Script and Chinese Civilization Inheritance and Development Project" under the project "Construction of Ancient Chinese Script Artificial Intelligence Recognition System" (Project No. G3829), and is supported by Graduate Work Department of Jilin University.

**Author contributions.** XY developed the research idea and provided valuable suggestions for this manuscript. YT conducted the experiments in this manuscript. YT, HP, and XF wrote this manuscript. CL provided specialized knowledge of ancient writing.

**Competing Interests.** The authors declare no competing interests.

## References

- [1] Epshtein, B., Ofek, E., Wexler, Y.: Detecting text in natural scenes with stroke width transform. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 2963–2970 (2010). IEEE
- [2] Lee, J.-J., Lee, P.-H., Lee, S.-W., Yuille, A., Koch, C.: Adaboost for text detection in natural scene. In: 2011 International Conference on Document Analysis and Recognition, pp. 429–434 (2011). IEEE
- [3] Jain, A.K., Yu, B.: Automatic text location in images and video frames. *Pattern recognition* **31**(12), 2055–2076 (1998)
- [4] Li, B., Dai, Q., Gao, F., Zhu, W., Li, Q., Liu, Y.: Hwobc-a handwriting oracle bone character recognition database. In: *Journal of Physics: Conference Series*, vol. 1651, p. 012050 (2020). IOP Publishing
- [5] Huang, S., Wang, H., Liu, Y., Shi, X., Jin, L.: Obc306: A large-scale oracle bone character recognition dataset. In: 2019 International Conference on Document Analysis and Recognition (ICDAR), pp. 681–688 (2019). <https://doi.org/10.1109/ICDAR.2019.00114>
- [6] Liu, G., Xing, J., Xiong, J.: Spatial pyramid block for oracle bone inscription detection. In: *Proceedings of the 2020 9th International Conference on Software and Computer Applications*, pp. 133–140 (2020)

- [7] Liu, G., Chen, S., Xiong, J., Jiao, Q.: An oracle bone inscription detector based on multi-scale gaussian kernels. *Applied Mathematics-a Journal of Chinese Universities Series B* **12**, 224–239 (2021)
- [8] Fu, X., Zhou, R.: Shape prior fusion for oracle bone inscriptions detection. In: *Proceedings of the 2024 7th International Conference on Image and Graphics Processing*, pp. 394–401 (2024)
- [9] Fu, X., Zhou, R., Yang, X., Li, C.: Detecting oracle bone inscriptions via pseudo-category labels. *Heritage Science* **12**(1), 107 (2024)
- [10] Bodla, N., Singh, B., Chellappa, R., Davis, L.S.: Soft-nms–improving object detection with one line of code. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5561–5569 (2017)
- [11] Ning, C., Zhou, H., Song, Y., Tang, J.: Inception single shot multibox detector for object detection. In: *2017 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, pp. 549–554 (2017). <https://doi.org/10.1109/ICMEW.2017.8026312>
- [12] Liao, M., Shi, B., Bai, X., Wang, X., Liu, W.: Textboxes: A fast text detector with a single deep neural network. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31 (2017)
- [13] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.-Y., Berg, A.C.: Ssd: Single shot multibox detector. In: *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pp. 21–37 (2016). Springer
- [14] Liao, M., Shi, B., Bai, X.: Textboxes++: A single-shot oriented scene text detector. *IEEE transactions on image processing* **27**(8), 3676–3690 (2018)
- [15] Zhou, X., Yao, C., Wen, H., Wang, Y., Zhou, S., He, W., Liang, J.: East: an efficient and accurate scene text detector. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5551–5560 (2017)
- [16] Tian, Z., Huang, W., He, T., He, P., Qiao, Y.: Detecting text in natural image with connectionist text proposal network. In: *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VIII 14*, pp. 56–72 (2016). Springer
- [17] Shi, B., Bai, X., Belongie, S.: Detecting oriented text in natural images by linking segments. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2550–2558 (2017)
- [18] Zhang, C., Liang, B., Huang, Z., En, M., Han, J., Ding, E., Ding, X.: Look more than once: An accurate detector for text of arbitrary shapes. In: *Proceedings of*

- the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10552–10561 (2019)
- [19] Liu, Y., He, T., Chen, H., Wang, X., Luo, C., Zhang, S., Shen, C., Jin, L.: Exploring the capacity of sequential-free box discretization network for omnidirectional scene text detection. *arXiv preprint arXiv:1912.09629* **3**, 15 (2019)
  - [20] Lyu, P., Liao, M., Yao, C., Wu, W., Bai, X.: Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes. In: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 67–83 (2018)
  - [21] Liao, M., Lyu, P., He, M., Yao, C., Wu, W., Bai, X.: Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes. *IEEE transactions on pattern analysis and machine intelligence* **43**(2), 532–548 (2021)
  - [22] Liao, M., Pang, G., Huang, J., Hassner, T., Bai, X.: Mask textspotter v3: Segmentation proposal network for robust scene text spotting. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pp. 706–722 (2020). Springer
  - [23] Baek, Y., Lee, B., Han, D., Yun, S., Lee, H.: Character region awareness for text detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9365–9374 (2019)
  - [24] Liao, M., Wan, Z., Yao, C., Chen, K., Bai, X.: Real-time scene text detection with differentiable binarization. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 11474–11481 (2020)
  - [25] He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9729–9738 (2020)
  - [26] Chen, X., Fan, H., Girshick, R., He, K.: Improved Baselines with Momentum Contrastive Learning (2020). <https://arxiv.org/abs/2003.04297>
  - [27] Chen, X., Xie, S., He, K.: An Empirical Study of Training Self-Supervised Vision Transformers (2021). <https://arxiv.org/abs/2104.02057>
  - [28] Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A Simple Framework for Contrastive Learning of Visual Representations (2020). <https://arxiv.org/abs/2002.05709>
  - [29] Chen, T., Kornblith, S., Swersky, K., Norouzi, M., Hinton, G.: Big Self-Supervised Models are Strong Semi-Supervised Learners (2020). <https://arxiv.org/abs/2006.10029>

- [30] Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., Joulin, A.: Unsupervised Learning of Visual Features by Contrasting Cluster Assignments (2021). <https://arxiv.org/abs/2006.09882>
- [31] Xie, J., Girshick, R., Farhadi, A.: Unsupervised Deep Embedding for Clustering Analysis (2016). <https://arxiv.org/abs/1511.06335>
- [32] Yang, J., Parikh, D., Batra, D.: Joint Unsupervised Learning of Deep Representations and Image Clusters (2016). <https://arxiv.org/abs/1604.03628>
- [33] Caron, M., Bojanowski, P., Joulin, A., Douze, M.: Deep clustering for unsupervised learning of visual features. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 132–149 (2018)
- [34] Caron, M., Bojanowski, P., Mairal, J., Joulin, A.: Unsupervised Pre-Training of Image Features on Non-Curated Data (2019). <https://arxiv.org/abs/1905.01278>
- [35] Feng, T., Wang, W., Wang, X., Yang, Y., Zheng, Q.: Clustering based point cloud representation learning for 3d analysis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 8283–8294 (2023)
- [36] Daneshfar, F., Saifee, B.S., Soleymanbaigi, S., Aeini, M.: Elastic deep multi-view autoencoder with diversity embedding. *Information Sciences* **689**, 121482 (2025)
- [37] Xing, J., Liu, G., Xiong, J.: Oracle bone inscription detection: a survey of oracle bone inscription detection based on deep learning algorithm. In: Proceedings of the International Conference on Artificial Intelligence, Information Processing and Cloud Computing, pp. 1–8 (2019)
- [38] Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence* **39**(6), 1137–1149 (2016)
- [39] Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: European Conference on Computer Vision, pp. 213–229 (2020). Springer
- [40] Sun, P., Zhang, R., Jiang, Y., Kong, T., Xu, C., Zhan, W., Tomizuka, M., Li, L., Yuan, Z., Wang, C., *et al.*: Sparse r-cnn: End-to-end object detection with learnable proposals. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14454–14463 (2021)
- [41] Li, B., Luo, D., Liang, Y., Yang, J., Ding, Z., Peng, X., Jiang, B., Han, S., Sui, D., Qin, P., *et al.*: Oracle bone inscriptions multi-modal dataset. arXiv preprint arXiv:2407.03900 (2024)
- [42] He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition

- (2015). <https://arxiv.org/abs/1512.03385>
- [43] Lin, T.-Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 936–944 (2017). <https://doi.org/10.1109/CVPR.2017.106>
- [44] Bottou, L., Curtis, F.E., Nocedal, J.: Optimization Methods for Large-Scale Machine Learning (2018). <https://arxiv.org/abs/1606.04838>
- [45] Loshchilov, I., Hutter, F.: Fixing weight decay regularization in adam. ArXiv [abs/1711.05101](https://arxiv.org/abs/1711.05101) (2017)
- [46] Powers, D.M.W.: Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation (2020). <https://arxiv.org/abs/2010.16061>
- [47] Zhang, H., Li, F., Liu, S., Zhang, L., Su, H., Zhu, J., Ni, L.M., Shum, H.-Y.: Dino: Detr with improved denoising anchor boxes for end-to-end object detection. arXiv preprint [arXiv:2203.03605](https://arxiv.org/abs/2203.03605) (2022)
- [48] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 10012–10022 (2021)
- [49] Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 779–788 (2016)
- [50] Maaten, L., Hinton, G.: Visualizing data using t-sne. *Journal of Machine Learning Research* **9**(86), 2579–2605 (2008)