# Contrast-weighted Dictionary Learning Based Saliency Detection for Remote Sensing Images

Zhou Huang[a], Huai-Xin Chen[a,*], Tao Zhou[b], Yun-Zhi Yang[c], Chang-Yin Wang[c], Bi-Yuan Liu[a]

[a]*University of Electronic Science and Technology of China, Chengdu, China*
[b]*Inception Institute of Artificial Intelligence (IIAI), Abu Dhabi, UAE*
[c]*CETC Special Mission Aircraft System Engineering Co.Ltd, Chengdu, China*

**Abstract**

Object detection is an important task in remote sensing image analysis. To reduce the computational complexity of redundant information and improve the efficiency of image processing, visual saliency models have been widely applied in this field. In this paper, a novel saliency detection model based on Contrast-weighted Dictionary Learning (CDL) is proposed for remote sensing images. Specifically, the proposed CDL learns salient and non-salient atoms from positive and negative samples to construct a discriminant dictionary, in which a contrast-weighted term is proposed to encourage the contrast-weighted patterns to be present in the learned salient dictionary while discouraging them from being present in the non-salient dictionary. Then, we measure the saliency by combining the coefficients of the sparse representation (SR) and reconstruction errors. Furthermore, by using the proposed joint saliency measure, a variety of saliency maps are generated based on the discriminant dictionary. Finally, a fusion method based on global gradient optimization is proposed to integrate multiple saliency maps. Experimental results on four datasets demonstrate that the proposed model outperforms other state-of-the-art methods.

*Keywords:* Contrast-weighted dictionary, Dictionary learning, Gradient

---

*Corresponding author.
*Email addresses:* `chowhuang23@gmail.com` (Zhou Huang), `huaixinchen@uesct.edu.cn` (Huai-Xin Chen ), `taozhou.ai@gmail.com` (Tao Zhou), `yangyz@cetca.net.cn` (Yun-Zhi Yang), `wangcy@cetca.net.cn` (Chang-Yin Wang), `lby9469@163.com` (Bi-Yuan Liu)

optimization, Remote sensing, Saliency detection

## 1. Introduction

Guided by our gaze, the human visual system (HVS) can quickly and automatically select regions of interest in complex scenes (known as the visual attention mechanism) [1]. This intelligent mechanism of the HVS has been extensively studied in the fields of psychology [2], neurobiology [3], and computer vision [4]. In the past two decades, research on visual saliency has advanced in two ways: eye fixation prediction in human vision [5] and salient object detection (SOD) in computer vision [6, 7]. The former focuses on predicting the eye fixations of an observer in a short time [8], whereas the latter aims to locate or segment the most prominent objects in a scene [9, 10, 11]. Because saliency detection can optimize the computing resources required for image analysis, visual saliency models are widely used in various fields of remote sensing (RS) image processing, including regional change detection [12], building detection [13] and oil tank detection [14].

The latest research [15] suggests that information is typically represented by a few simultaneously active neurons. Importantly, while the retina receives a lot information, only a small amount of useful data is transmitted to nerve cells in the visual cortex for processing. This representation of information is known as a sparse representation (SR) [16]. The principle of SR is to represent the signal by a linear combination of a series of base vectors in the over-complete dictionary, and that linear combination must be sparse [17]. In recent years, image structure analysis based on SR has been widely used in computer vision and image processing. At the same time, SR theory has been introduced into the field of image saliency detection [18, 19]. However, there are two key problems with SR-based SOD methods: the construction of the SR dictionary and the criteria for saliency measure.

In the construction of dictionaries, most of the early methods used independent component analysis (ICA) to sample numerous image patches from

various kinds of natural images to generate basic atoms [8]. However, these basic atoms cannot create a perfect SR of the detection image without information loss because some features of the training image cannot be accurately captured by the predetermined basic atom. Other SR-based methods [20, 21] usually use the areas around the detected patches for dictionary construction. However, as [22] showed, when the salient object has a high contrast with the surrounding patches, such methods usually assign higher values to the edges of the salient object rather than the entire object. In addition, in [23], a multi-view joint SR framework that simultaneously considers the inherent contextual structures among instances improved the performance and robustness of the learned dictionary. Recently, the background prior [24] was introduced into SR-based saliency detection methods, which assumes that non-salient parts of the image are usually distributed on the boundary. Under this assumption, patches or superpixels near the boundaries of the image are usually selected to build the background dictionary [25, 26]. However, when the salient object is near the image boundary, some foreground regions are included in the background dictionary, which causes them to be mistakenly detected as background regions. Also, if the background regions near the boundary of the image have distinct features, some background regions will be incorrectly marked as foreground. Moreover, the training sample patches usually have their own characteristic features, such as intensity and contrast, but these are usually disregarded in most existing SOD methods, resulting in salient objects in a scene with similar background and foregrounds being unevenly highlighted.

As for saliency measurement criteria, saliency detection methods based on SR define this in terms of reconstruction error or sparsity of representation coefficients (that is, using the $l_0$-norm to calculate the coding length) [18, 19, 27]. These methods also usually add sparse constraints to sparse coefficients to achieve sparse coding of image patches, and they calculate the saliency of image patches by minimizing the sum of the reconstruction errors. Therefore, these representation methods are more sensitive to non-Gaussian noise rather than outliers representing coefficients.

Through our research, we have found that two or more temporary saliency maps are generated in most saliency detection models. Among these methods, some determine the fusion weights through simple weights [19, 27] or experimental effects [8]. Other methods determine the optimal image fusion weights through methods such as least-squares estimation of training data [28] or Bayesian inference [18], but do not consider the connections between multiple saliency maps.

To solve the above problems, we propose an SR method based on Contrast-weighted Dictionary Learning (CDL) for saliency detection. Specifically, this paper uses the positive and negative samples generated by the salient and non-salient regions in the image as a template for dictionary learning. Inspired by the online dictionary learning algorithm [17], to solve the problem of dictionary learning we also propose an online discriminant CDL algorithm, which effectively overcomes the shortcomings of some methods using background priors. To determine saliency, we use the $l_2$-norm to measure the sparsity of sparse coefficients, combined with the $l_{1,2}$-norm to calculate the sparse reconstruction errors and improve the expression of outliers in the sparse coefficient. For the various saliency maps generated by calculating representation coefficients, we propose an image fusion method based on global gradient optimization to integrate multiple salient images. To summarize, the main contributions of this paper are as follows:

(1) Considering the features of the training sample patch itself, we propose a novel atomic learning formula based on contrast weights. Further, we use an online discriminative CDL to solve the formula.

(2) We use the $l_2$-norm to measure the sparsity of sparse coefficients, use the $l_{2,1}$-norm to measure the sparsity of the reconstruction errors, and then combine the two measures to improve the expression of outliers in the representation coefficients.

(3) We use a salient map fusion method based on global gradient optimization to integrate multiple saliency maps. This method optimizes the image

4

fusion effect by establishing the relation between saliency maps.

The rest of this paper is organized as follows: Section 2 briefly reviews related work. Section 3 describes the SOD method in detail. Sections 4 and 5 give the experimental analysis and conclusions.

## 2. Related Work

In recent years, more and more researchers are committed to the work of SOD [29, 30]. Several review papers [1, 31] have investigated and discussed many of the most advanced SOD methods in detail. In this section, we review the work most relevant to ours, including SOD based on sparse representations and the application of saliency detection in optical RS images. Further, SOD based on deep learning is another hot topic in recent years and will be briefly reviewed in this section.

### 2.1. SOD based on sparse representations

In recent years, SR theory has been gradually addressed in the field of saliency detection. Generally, SR based saliency detection methods need to first construct an over-complete dictionary, then sparsely represent an input image through the dictionary, and finally measure the saliency according to the SR coefficients or reconstruction errors. In [8], the construction of the dictionary was learned by applying ICA on the image patches sampled from each position of the input image and using the reconstruction errors to measure the saliency. In the method of [20] the image patches around the central patch were used for SR, and the saliency was measured by the coding length or residual. These methods usually give higher saliency values to the object boundaries, because both the background and foreground are included in the dictionary.

Later, the background prior method [24] was proposed. As an extension of this, some methods then [25, 26] used patches or superpixels near the image boundary as background templates to construct a global background dictionary and sparsely reconstruct the image. Recently, in [32], a SOD method was

proposed based on two-stage graphs, taking into account the consistency of adjacent spaces between graph nodes and the consistency of regional spaces, while improving the accuracy of SOD in complex scenes.

## 2.2. Application of SOD in RS images

Due to the rapid development of massive RS data and the complexity of RS scenes, many traditional methods of processing natural images are not suitable for RS images. As one method of data compression and rapid screening, saliency detection can effectively process RS data. Importantly, there are several essential similarities between SOD/ target detection and extraction in RS images. For instance, both extract regions of interest in an image based on the saliency of a particular task or target. As image processing and RS technology have developed, saliency detection has been widely used in the field of RS. Many researchers have combined visual saliency and image interpretation to accomplish specific target detection, such as regional change detection [12], airport detection [33], building detection [13] and oil tank detection [14]. For example, Yao et al. [33] proposed a coarse-to-fine airport saliency detection model. At the coarse layer, combined with contrast and linear density clues, a goal-oriented saliency model was established to quickly locate airport candidate regions. Later, Li [13] et al. proposed a two-step building extraction method based on saliency cues, designed a saliency estimation algorithm for building objects, extracted saliency cues in a local region of each candidate building, and integrated them into a probability model to get the final building extraction results. However, these methods do not involve road detection methods based on saliency in RS images.

## 2.3. SOD based on deep learning

Recently, SOD methods based on deep learning have attracted more attention. Zhang et al. [34] proposed a SOD model based on fully a convolutional neural network by introducing a gated two-way message passing module to integrate multi-level features. In [35], a predict-refine architecture and a new

hybrid loss for boundary-aware SOD were proposed to pay more attention to the boundary quality of salient objects. In another work focusing on salient edge information [36], an essential pyramid attention structure for SOD was designed to enhance the representation ability of the corresponding network layer.

In order to prevent SOD methods for RGB images from failing [24, 31, 37] when processing complex scenes, as is common in recent saliency detection works dedicated to RGB-D, [38] introduced the probabilistic RGB-D saliency detection network via conditional variational autoencoders to model human annotation uncertainty and generate multiple saliency maps for each input image by sampling in the latent space. Further, Fan et al. [39] constructed a 1K high-resolution saliency person dataset, and proposed a baseline architecture called the Deep Depth-Depurator Network for saliency detection. In addition, binocular stereo cameras are widely used in various tasks of RS photogrammetry, which makes it possible to use supplementary depth information to further accurately detect and identify targets in the field of RS.

## 3. Proposed Saliency Detection Model

This section describes the proposed saliency detection model in detail. As shown in Fig. **??**, the model includes three main parts: CDL-based discriminant dictionary learning, saliency maps generation and fusion.

### 3.1. Contrast-weighted dictionary learning formula

In the image processing method based on SR, an image patch is usually represented by a linear combination of a few atoms in an over-complete dictionary $\mathbf{D} = \{\mathbf{d}_i\}_{i=1}^{k} \in \mathbb{R}^{n \times k}$; that is, the image patch $\mathbf{x} \in \mathbb{R}^n$ is estimated by dictionary $\mathbf{D}$ and the calculated sparse coefficients $\boldsymbol{\alpha} \in \mathbb{R}^k$ . The equation is

$$\mathbf{x} = \mathbf{D}\boldsymbol{\alpha} \quad s.t. \quad \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2 \leq \xi, \tag{1}$$

where $\|\bullet\|_2$ is the $l_2$-norm used to measure the deviation and $\xi$ is the error. Within the feasible set, the solution that minimizes the number of nonzero
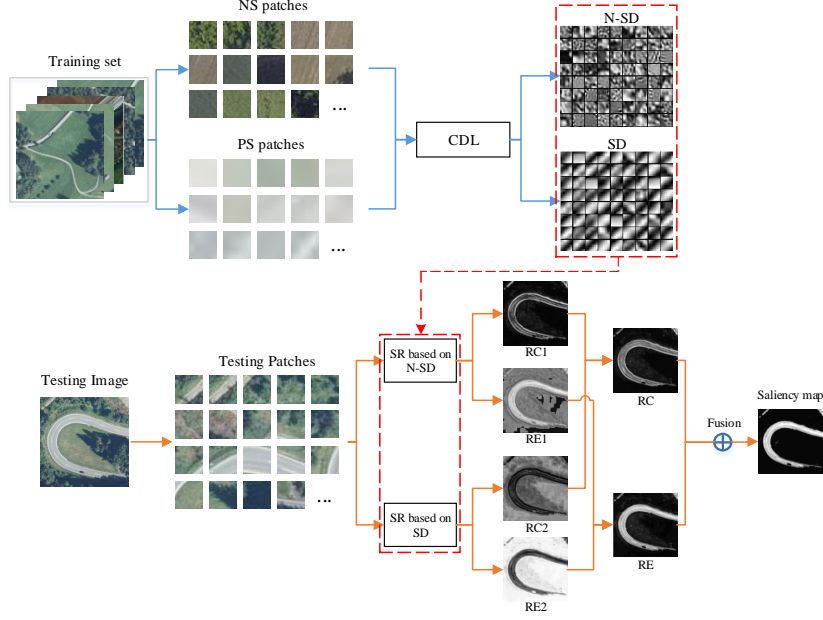
7

Figure 1: Diagram of the proposed SOD method. NS: negative sample, PS: positive sample, CDL: contrast-weighted dictionary learning, N-SD: non-salient dictionary, SD: salient dictionary, SR: sparse representation, RC: representation coefficient, RE: reconstruction error.

sparse coefficients is undoubtedly an attractive representation. This form can be expressed as:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^k} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \quad s.t. \quad \|\boldsymbol{\alpha}\|_0 \leq L, \tag{2}$$

where $L$ is the sparsity of the coefficients $\boldsymbol{\alpha}$. In Eq. 2, the atoms in $\mathbf{D}$ represent the smallest unit in the reconstructed image patches. Here, the atoms in $\mathbf{D}$ need to be learned from the training patches $\mathbf{X} = \{\mathbf{x}\}_{i=1}^m$, which can be achieved by [19]

$$\min_{\mathbf{D}, \mathbf{A}} \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1 \right), \tag{3}$$

where $\lambda$ is the trade-off between the reconstruction errors $\|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2$ and the sparsity of the coefficient $\|\boldsymbol{\alpha}_i\|_1$ and $\mathbf{A} = \{\boldsymbol{\alpha_i}\}_{i=1}^m$ is the SR coefficients set corresponding to $\mathbf{X}$. According to Eq. 3, we study salient and non-salient dictionary learning based on contrast-weighted atoms. As one of the features of an image, contrast plays an important role in both local and global saliency

detection methods [40, 41]. In order to enhance the ability of the base atom to learn the image contrast features and improve the sensitivity to the contrast of the surrounding pixels, a novel contrast-weighted term is incorporated in our formulation to encourage/discourage the contrast-weighted patterns in the learned salient/non-salient dictionary, respectively. More specifically, in our weight function, the weight of each pixel in the base atom is calculated by the relative brightness contrast of the pixels and the corresponding training sample patch $\mathbf{x}_i$. Thus, the weight of the $j$-th pixel $p_{ij}$ in the $i$-th sample patch is:

$$\mathbf{w}\left(p_{ij}\right) = \frac{Lum\left(p_{ij}\right) - \text{mean}\left(Lum\left(\mathbf{x}_i\right)\right)}{\max\left(Lum\left(\mathbf{x}_i\right)\right)}, \tag{4}$$

where $Lum\left(\cdot\right)$ is the luminance value operator for calculating the sample patch, and mean $\left(\cdot\right)$ and max $\left(\cdot\right)$ are the average value operator and the maximum value operator, respectively. Note that in the actual calculation, the $i$-th sample patch is treated as a column vector; that is, $\mathbf{w}\left(p_{ij}\right)$ can be expressed as $\mathbf{w}_{ij}^T \in \mathbb{R}^{1 \times n}$, and $n$ is the number of pixels of the sample patch. Upon $\mathbf{W}_i^T \in \left\{\mathbf{w}_{ij}^T\right\}_{i,j=1}^{m,n}$, the contrast weight term can be designed by $\left\|\mathbf{W}_i^T \mathbf{D}\right\|_2^2$, which quantifies the degree of weighted contrast. Thus, given the contrast weight term, by rewriting Eq. 3, we have the following formula for optimizing the salient and non-salient dictionary learning:

$$\min_{\mathbf{D}^H, \mathbf{A}^H} \frac{1}{m^H} \sum_{i=1}^{m^H} \left(\frac{1}{2}\left\|\mathbf{x}_i^H - \mathbf{D}^H \boldsymbol{\alpha}_i^H\right\|_2^2 + \lambda_1 \left\|\boldsymbol{\alpha}_i^H\right\|_1 + \lambda_2 \left\|\mathbf{W}_i^T \mathbf{D}^H\right\|_2^2\right), H = \left\{\text{ P} \vee \text{N}\right\}, \tag{5}$$

where $m^H$ represents the number of positive or negative samples, $\left\{\boldsymbol{\alpha}_i^H\right\}_{i=1}^{m^H}$ is the SR coefficient of positive or negative sample patches, and $\mathbf{D}^H$ is a salient or non-salient dictionary trained from positive or negative samples. The meaning of $\lambda_1$ is the same as that of $\lambda$ in formula (3), and $\lambda_2$ (a very small positive number) is a regularization parameter that controls the influence of contrast-weighted terms.

### 3.2. The solution to the dictionary learning formulation

We can learn the salient and non-salient dictionaries through (5). Because the online dictionary learning algorithm [17] can deal with large dynamic datasets and is faster than the batch algorithm, we propose the CDL algorithm to solve (5). Similar to the standard dictionary learning algorithm, we divide the optimization problem in (5) into two subprocesses to solve alternately; namely, the SR and dictionary update. Specifically, the initialization training dictionary $\mathbf{D}^{\mathrm{H}}$ is generally obtained by randomly sampling the training sample set. Thus, the first step is to fix $\mathbf{D}^{\mathrm{H}}$ , and the sparse coefficient $\mathbf{A}^{\mathrm{H}} = \left\{ \boldsymbol{\alpha}_i^{\mathrm{H}} \right\}_{i=1}^{m^{\mathrm{H}}}$ can then be obtained by the SR method. The second step is to fix $\mathbf{A}^{\mathrm{H}}$ , and the updated dictionary $\mathbf{D}^{\mathrm{H}}$ can then be solved by the dictionary update method. The first and second steps of the iteration are done until convergence is reached.

### 3.2.1. Sparse representation

From the above, it can be seen that the solution to (5) is an iterative optimization process, assuming that in the $i$-th iteration, $\mathbf{x}_t^{\mathrm{H}}$ is a randomly selected image patch from the training set, and $\boldsymbol{\alpha}_t^{\mathrm{H}}$ is the coefficient of $\mathbf{x}_t^{\mathrm{H}}$ obtained by the $(t-1)$-th updated dictionary $\mathbf{D}_{t-1}^{\mathrm{H}}$ through the SR algorithm. Because the contrast-weighted term $\lambda_2 \left\| \mathbf{W}_i^T \mathbf{D}^{\mathrm{H}} \right\|_2^2$ in (5) is independent of the sparse coefficient $\boldsymbol{\alpha}_i^{\mathrm{H}}$ , the sparse coefficient a in the $i$-th iteration can be expressed as

$$\alpha_t^{\mathrm{H}} \triangleq \underset{\alpha_t^{\mathrm{H}} \in \mathbb{R}^k}{\arg\min} \frac{1}{2} \left\| \mathbf{x}_t^{\mathrm{H}} - \mathbf{D}_{t-1}^{\mathrm{H}} \boldsymbol{\alpha}_t^{\mathrm{H}} \right\|_2^2 + \lambda_1 \left\| \boldsymbol{\alpha}_t^{\mathrm{H}} \right\|_1, \tag{6}$$

The SR problem of the above fixed dictionary is the $l_1$-regularized linear least square problem. In this paper, the LARS-Lasso algorithm [42] is used to solve this problem.

### 3.2.2. Dictionary update

After the SR step of the $t$-th iteration, the sparse coefficient $\left\{ \alpha_i^{\mathrm{H}} \right\}_{i=1}^{t}$ of the image patch $\left\{ \mathbf{x}_i^{\mathrm{H}} \right\}_{i=1}^{t}$ after training is obtained. In the $t$-th iteration, with fixed $\alpha_i^{\mathrm{H}}$, the dictionary can be updated using the following optimization function

according to (5):

$$\mathbf{D}_t^{\mathrm{H}} \triangleq \underset{\mathbf{D}_t^{\mathrm{H}} \in \mathbb{R}^{n \times k}}{\arg \min} \frac{1}{t} \sum_{i=1}^{t} \left( \frac{1}{2} \left\| \mathbf{x}_i^{\mathrm{H}} - \mathbf{D}_{t-1}^{\mathrm{H}} \boldsymbol{\alpha}_t^{\mathrm{H}} \right\|_2^2 + \lambda_1 \left\| \boldsymbol{\alpha}_t^{\mathrm{H}} \right\|_1 + \lambda_2 \left\| \mathbf{W}^T \mathbf{D}_{t-1}^{\mathrm{H}} \right\|_2^2 \right), \tag{7}$$

where $\mathbf{D}_t^{\mathrm{H}}$ is the discriminant dictionary obtained after the $t$-th iterative learning.

Because the patch coordinate descent algorithm [43] has the advantages of no parameters and no need for any learning rate adjustment, we updated each atom of the dictionary using this algorithm. For example, the $j$-th atom $\mathbf{d}_{j,t}^{\mathrm{H}}$ for updating the dictionary in the $t$-th iteration is calculated by

$$\mathbf{d}_{j,t}^{\mathrm{H}} = \mathbf{d}_{j,t-1}^{\mathrm{H}} - \frac{\sigma}{t} \frac{\partial}{\partial \mathbf{d}_j^{\mathrm{H}}} \left[ \sum_{i=1}^{t} \left( \frac{1}{2} \left\| \mathbf{x}_i^{\mathrm{H}} - \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \boldsymbol{\alpha}_i^{\mathrm{H}} \right\|_2^2 + \lambda_1 \left\| \boldsymbol{\alpha}_i^{\mathrm{H}} \right\|_1 \right. \right.$$
$$\left. \left. + \lambda_2 \left\| \mathbf{W}_j^T \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \right\|_2^2 \right) \right]_{|\mathbf{d}_{j,t-1}^{\mathrm{H}}} . \tag{8}$$

For convenience, let

$$\mathbf{M} = \sum_{i=1}^{t} \left( \frac{1}{2} \left\| \mathbf{x}_i^{\mathrm{H}} - \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \boldsymbol{\alpha}_i^{\mathrm{H}} \right\|_2^2 + \lambda_1 \left\| \boldsymbol{\alpha}_i^{\mathrm{H}} \right\|_1 + \lambda_2 \left\| \mathbf{W}_j^T \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \right\|_2^2 \right). \tag{9}$$

In Eq. 8, $\sigma$ is the learning rate of gradient descent, and $\widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} = [\mathbf{d}_{1,t}^{\mathrm{H}}, \mathbf{d}_{2,t}^{\mathrm{H}}, \cdots,$ $\mathbf{d}_{j,t}^{\mathrm{H}}, \mathbf{d}_{j+1,t-1}^{\mathrm{H}}, \cdots, \mathbf{d}_{k,t-1}^{\mathrm{H}}]$, Note that the only variable that must be updated is $\mathbf{d}_{j,t}^{\mathrm{H}}$ in $\widehat{\mathbf{D}}_{j,t}^{\mathrm{H}}$. After Eq. 8 for the current iteration, the previous $j$ atoms, that is, $\left\{ \mathbf{d}_{1,t}^{\mathrm{H}}, \mathbf{d}_{2,t}^{\mathrm{H}}, \cdots, \mathbf{d}_{j,t}^{\mathrm{H}} \right\}$, are updated. Using the trace $Tr(\bullet)$ of the matrix to represent the $l_2$-norm and then expressing it as the derivative of $\mathbf{d}_j^{\mathrm{H}}$, Eq. 9 can be rewritten as

$$\frac{\partial}{\partial \mathbf{d}_j^{\mathrm{H}}} (\mathbf{M})_{|\mathbf{d}_{j,t-1}^{\mathrm{H}}} = \frac{1}{2} \frac{\partial}{\partial \mathbf{d}_j^{\mathrm{H}}} Tr \left[ \left( \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \right)^T \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \mathbf{B}_t^{\mathrm{H}} \right] - \frac{\partial}{\partial \mathbf{d}_j^{\mathrm{H}}} Tr \left[ \left( \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \right)^T \mathbf{C}_t^{\mathrm{H}} \right]$$
$$+ \frac{\partial}{\partial \mathbf{d}_j^{\mathrm{H}}} Tr \left[ \lambda_2 t \mathbf{W}_j^T \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \left( \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \right)^T \mathbf{W}_j \right], \tag{10}$$

where $\mathbf{B}_t^{\mathrm{H}}$ and $\mathbf{C}_t^{\mathrm{H}}$ are defined as $\sum_{i=1}^{t} \alpha_i^{\mathrm{H}} \left( \alpha_i^{\mathrm{H}} \right)^T$ and $\sum_{i=1}^{t} \mathbf{x}^{\mathrm{H}} \left( \alpha_i^{\mathrm{H}} \right)^T$, which

11

Table 1: Summary of the CDL algorithm

---

**Algorithm 1**. Online discriminant dictionary learning algorithm based on weighted contrast

---

**Input:** Vectorised training patches $\mathbf{X}^{\mathrm{H}} \in \mathbb{R}^{n \times m}$.
**Output:** The learned dictionary $\mathbf{D}^{\mathrm{H}} \in \mathbb{R}^{n \times k}$.
**Initialization:** The contrast-weighted matrix $\mathbf{W}^{T}$ is obtained by Eq. 4;
                     Randomly select the samples in the training set to fill $\mathbf{D}_0^{\mathrm{H}}$;
                     Set $\mathbf{B}_0^{\mathrm{H}} \in \mathbb{R}^{k \times k}$ and $\mathbf{C}_0^{\mathrm{H}} \in \mathbb{R}^{n \times k}$ to zero matrices;
                     Regularization parameter $\lambda_1$ and $\lambda_2$ ;
                     Number of iterations $T$.
1. **For** $t = 1$ to $T$ **do**
2.        Randomly select the image patches $\mathbf{X}^{\mathrm{H}}$ from the training set $\mathbf{x}_t^{\mathrm{H}} \in \mathbb{R}^{k \times 1}$.
3.        Sparse coding:
                 Obtained $\alpha_t^{\mathrm{H}} \in \mathbb{R}^{k \times 1}$ by solving Eq. 6 with LARS-Lasso [60] algorithm.
4.        Update $\mathbf{B}_t^{\mathbf{H}}$ and $\mathbf{C}_t^{\mathbf{H}}$:
                 $\mathbf{B}_t^{\mathbf{H}} = \sum_{i=1}^{t} \alpha_i^{\mathrm{H}} \left( \alpha_i^{\mathrm{H}} \right)^{T} = \mathbf{B}_{t-1}^{\mathbf{H}} + \alpha_t^{\mathrm{H}} \left( \alpha_t^{\mathrm{H}} \right)^{T}$,
                 $\mathbf{C}_t^{\mathbf{H}} = \sum_{i=1}^{t} \mathbf{x}_i^{\mathrm{H}} \left( \alpha_i^{\mathrm{H}} \right)^{T} = \mathbf{C}_{t-1}^{\mathbf{H}} + \mathbf{x}^{\mathrm{H}} \left( \alpha_t^{\mathrm{H}} \right)^{T}$.
5.        Dictionary update:
                 **For** $t = 1$ to $T$ **do**
                 $\mathbf{d}_{j,t}^{\mathrm{H}} = \mathbf{d}_{j,t-1}^{\mathrm{H}} - \frac{1}{\mathbf{B}_j^{\mathrm{H}}(j,j)} \left( \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \mathbf{b}_{j,t}^{\mathrm{H}} - \mathbf{c}_{j,t}^{\mathrm{H}} \right) - 2\lambda_2 \sigma \mathbf{W}_j \mathbf{W}_j^{T} \mathbf{d}_{j,t-1}^{\mathrm{H}}$,
                 **End For**
6.       Obtain the discriminant dictionary $\mathbf{D}_t^{\mathrm{H}} = \left[ \mathbf{d}_{1,t}^{\mathrm{H}}, \mathbf{d}_{2,t}^{\mathrm{H}}, \cdots, \mathbf{d}_{k,t}^{\mathrm{H}} \right]$ for the current iteration.
7. **End For**
8. Return: The learned dictionary $\mathbf{D}^{\mathrm{H}} = \mathbf{D}_T^{\mathrm{H}}$.

---

refer to storing all the information of the sparse coefficients and sparsely represented image patches of all previous iterations, respectively. According to the derivative calculation rule of the matrix trace, Eq. 10 can be expressed as

$$\frac{\partial}{\partial \mathbf{d}_j^{\mathrm{H}}} \left( \mathbf{M} \right)_{|\mathrm{d}_{j,t-1}^{\mathrm{H}}} = \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \mathbf{b}_{j,t}^{\mathrm{H}} - \mathbf{c}_{j,t}^{\mathrm{H}} + 2\lambda_2 t \mathbf{W}_j \mathbf{W}_j^{T} \mathbf{d}_j^{\mathrm{H}}, \tag{11}$$

where $\mathbf{b}_{j,t}^{\mathrm{H}}$ and $\mathbf{c}_{j,t}^{\mathrm{H}}$ represent the $j$-th columns of $\mathbf{B}_{j,t}^{\mathrm{H}}$ and $\mathbf{C}_{j,t}^{\mathrm{H}}$ respectively. Thus Eq. 8 can be rewritten as

$$\mathbf{d}_{j,t}^{\mathrm{H}} = \mathbf{d}_{j,t-1}^{\mathrm{H}} - \frac{\sigma}{t} \left( \widehat{\mathbf{D}}_{j,t}^{\mathrm{H}} \mathbf{b}_{j,t}^{\mathrm{H}} - \mathbf{c}_{j,t}^{\mathrm{H}} \right) - 2\lambda_2 \sigma \mathbf{W}_j \mathbf{W}_j^{T} \mathbf{d}_{j,t-1}^{\mathrm{H}}. \tag{12}$$

According to [17], the $\sigma/t$ in Eq. 12 can be expressed approximately as $1/\mathbf{B}_j^{\mathrm{H}}(j,j)$ . When all the atoms $\left\{ \mathbf{d}_{j,t}^{\mathrm{H}} \right\}_{j=1}^{k}$ are updated, the dictionary $\mathbf{D}_t^{\mathrm{H}}$ completes the $t$-th learning.

In summary, after the iterative SR and dictionary update steps, we obtain salient and non-salient dictionaries. Table 1 summarizes our CDL algorithm.

### 3.3. Saliency image generation

This subsection describes the saliency measurement criteria based on SR coefficients and reconstruction errors.

### 3.3.1. Saliency measure based on sparse representation coefficients

In the saliency detection process, the saliency of each pixel can be measured to a certain extent by the representation coefficient of an image patch centered on the pixel, where the different representation coefficients of the image patch $\mathbf{x}_i$ are calculated by the discrimination dictionary $\mathbf{D}^H$ through the following formula:

$$\alpha_i^H = \arg\min \frac{1}{2} \left\| \mathbf{x}_i - \mathbf{D}^H \alpha_i^H \right\|_2^2 + \lambda_1 \left\| \alpha_i^H \right\|_1, \tag{13}$$

where $\lambda_1$ has same meaning as $\lambda$ in Eq. 3. As shown in Fig. 2, when using the salient dictionary for sparse reconstruction, non-salient image patches obtain their SR coefficients with high energy, while salient image patches obtain their SR coefficients with lower energy. This is because the salient dictionary has high contrast with non-salient image patches, and the saliency image patches have low contrast. On the basis of this observation, we define the saliency measure of a pixel as:

$$S_{\mathbf{A}(i)} = 1 - \exp\left( -\frac{\left\| \alpha_i^N \right\|_2^2 - \left\| \alpha_i^P \right\|_2^2}{2\eta_{\mathbf{A}}^2} \right), \tag{14}$$

where $\alpha_i^N$ and $\alpha_i^P$ represent the representation coefficients obtained by Eq. 13 for the image patch centered on pixel $i$, and $\eta_{\mathbf{A}}$ is a scalar parameter, which is set to 1 in the experiment.

### 3.3.2. Saliency measurement based on reconstruction error

Reconstruction error is widely used in saliency detection based on SR. Generally speaking, an image patch has a larger relative reconstruction error for the discriminant dictionary, so it will have a greater saliency value. Therefore, we
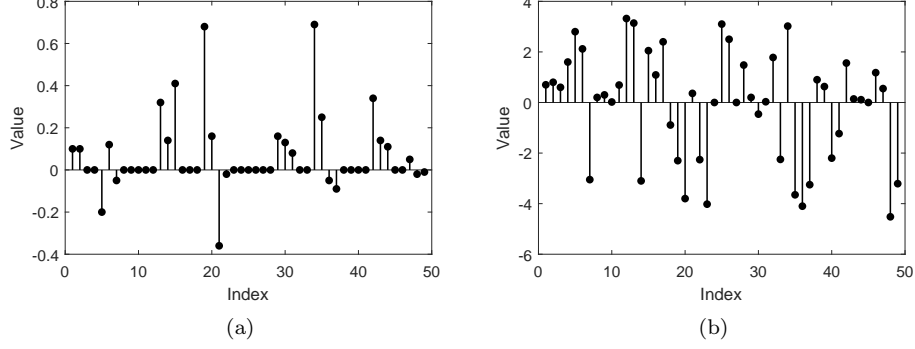
13

Figure 2: Comparison of SR coefficients using salient dictionary. (a) SR coefficients of salient patches. (b) SR coefficients of non-salient patches.

define the saliency measure of pixels based on SR coefficients as:

$$
S_{\mathbf{R}(i)} = 1 - \exp\left(-\frac{\min\limits_{\alpha_i^N}\left\|x_i - D^N\alpha_i^N\right\|_{2,1} - \min\limits_{\alpha_i^P}\left\|x_i - D^P\alpha_i^P\right\|_{2,1}}{2\eta_{\mathbf{R}}^2}\right), \quad (15)
$$

where $x_i$ is the image patch centered on a pixel $i$ , $D^N$ and $D^P$ represent the non-salient and salient dictionary, respectively, $\alpha_i^N$ and $\alpha_i^P$ are the representation coefficients obtained by the discriminant dictionary, and $\eta_{\mathbf{R}}$ is the scale parameter and is set to 1 in the experiment.

### 3.4. Saliency map fusion

In the field of information fusion, information fusion methods can achieve better results than a single information source as long as there are appropriate fusion criteria. The traditional pixel-level saliency map fusion method generates a fused image through the weighted sum of multiple saliency maps, which can be expressed as:

$$
S_{fused}(x, y) = \sum_{n=1}^{N} W_n(x, y)\, S_n(x, y), \quad (16)
$$

14

where $N$ is the number of saliency maps to be fused, $S_n(x, y)$ is the pixel intensity of the $n$-th saliency map at $(x, y)$, and $W_n(x, y)$ is the weight of the importance of pixel $S_n(x, y)$ at $(x, y)$. Therefore, the key to fusion is designing a reasonable weight.

With this in mind, and based on the observation of the cumulative histogram of pixel intensity (the histogram integral along the pixel intensity axis as shown in Fig. 3), we propose a weight function to suppress the background region and highlight the foreground region in the fusion of saliency maps.

Fig. 3 shows an example of a cumulative histogram of a coefficient representation map, a reconstruction error map to be fused, and an optimized fusion map. The cumulative histogram in the optimized fusion map increases sharply at the beginning; in other words, there are a significant number of pixels in this interval, and the intensity of the surrounding pixels has a small change with a larger gradient than that of the saliency image to be fused. In the middle region of pixel intensity (0.2 to 0.8), the cumulative histogram changes slowly, indicating that there are relatively few pixels in this region, and the surrounding pixel intensity has a greater change that has a smaller gradient compared with the saliency maps to be fused. The analysis for the interval where the pixel intensity is close to 1 follows the same rule as above. Therefore, when the pixels are in the range of a cumulative histogram with a large gradient, they need to be given a higher weight during image fusion. Formally, we can express this
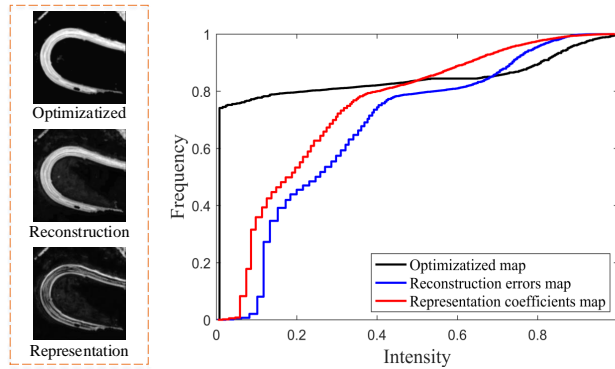


Figure 3: Cumulative histogram comparison of different saliency maps.

observation as:

$$W_n\left(x,y\right) = \frac{Grad_n\left(I_n\left(x,y\right)\right)}{\sum_{n=1}^{N} Grad_n\left(I_n\left(x,y\right)\right) + \varphi}, \qquad (17)$$

where $\varphi$ prevents the occurrence of a very small positive number with zero denominators, and $Grad_n\left(I_n\left(x,y\right)\right)$ is the gradient of the cumulative histogram at pixel intensity $I_n\left(x,y\right)$. Because the cumulative histogram is the statistical information of all pixels, the gradient in Eq. 17 is not the local gradient around the pixel; we call it the global gradient. Using the weights obtained above, we can fuse several saliency maps obtained by representation coefficients and reconstruction errors according to Eq. 16.

## 4. Experiments

In this section, we first introduce the constructed RS image dataset and three other popular datasets containing natural images, and then explain the dictionary training strategy, evaluation metrics and implementation details. Finally, we compare the proposed method with nine state-of-the-art methods.

### 4.1. Experimental setup

#### 4.1.1. Datasets

To the best of our knowledge, no publicly available dataset of optical RS images can be used for road detection. Therefore, we collected 300 optical RS images to build a dataset for road saliency detection, which we called "2RSOD", and manually annotated each image, pixel-wise. Most of the original optical RS images were collected from Google Earth, and the rest were collected from existing optical RS image datasets, including DOTA [44] and NWPU VHR-10 [45]. This 2RSOD dataset is challenging because the spatial resolutions of the images are diverse, including $300 \times 300$ , $500 \times 500$ and $1024 \times 1024$. Further, image backgrounds tend to be complicated and cluttered, often including buildings, trees, rivers, and shadows. The sizes, numbers, and shapes of the salient objects also vary. Some sample images from the constructed 2RSOD dataset are shown

16

Table 2: Summary of evaluation metrics

| Metric | Mathematical Expression |
|---|---|
| Precision-Recall $(PR)\uparrow$ | $Precision\,(P): \frac{\lvert S \cap G\rvert}{\lvert S\rvert},\ Recll\,(R): \frac{\lvert S \cap G\rvert}{\lvert G\rvert}$ |
| F-measure$(F_\beta)\uparrow$ | $F_\beta = \left(1+\beta^2\right)\frac{P*R}{\beta^2 P + R}, \beta^2 = 0.3$ |
| S-measure$(S_\alpha)\uparrow$ | $S = \alpha * S_0 + (1-\alpha) * S_r, \alpha = 0.5$ |
| E-measure$(E_\xi)\uparrow$ | $E = \frac{1}{W*H}\sum_{i=1}^{W}\sum_{i=1}^{H}\phi_{FM}\,(i,j)$ |
| Adaptive threshold $(F_{adp})\uparrow$ | $Thr = \frac{2}{W*H}\sum_{i=1}^{W}\sum_{i=1}^{H}S\,(i,j)$ |
| Mean absolute error $(MAE\ \mathcal{M})\downarrow$ | $MEA = \frac{1}{W*H}\sum_{i=1}^{W}\sum_{i=1}^{H}\lvert S\,(i,j) - G\,(i,j)\rvert$ |

Note: $\uparrow$ & $\downarrow$ denote larger and smaller is better.
$S$:saliency image $\qquad\qquad\qquad\qquad\qquad$ $S_o$:target perception structure [49]
$G$:corresponding annotation map $\qquad\quad$ $\phi$:enhanced contrast matrix [50]
$\lvert\cdot\rvert$:calculates the number of nonzero entries $\quad$ $W$:width of the image
$S_r$:similarity measurement of the region[49] $\quad$ $H$:height of the image

in Fig. 4. In addition, we evaluate CDL on three other benchmark natural image datasets, including ECSSD [46] with 1000 images, PASCAL-S [47] with 850 images, and DUT-OMRON [48] with 5168 images.



Figure 4: Sample images from the constructed 2RSOD dataset. The first row shows the optical RS images. The second row provides the pixel-wise annotations.
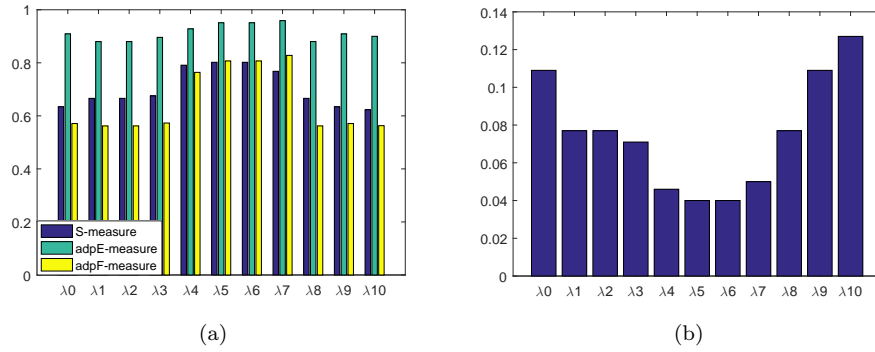
### 4.1.2. Evaluation metrics

To quantitatively evaluate the performance of various methods, we adopt six evaluation metrics. Table 2 summarizes these metrics.

### 4.1.3. Parameter settings

All parameter settings related to our experiment are summarized in Table 3. For 2RSOD and the other three natural image datasets, we select 240 images from each as the training sets for the discriminant dictionaries, and the remaining images as the test sets. For the dictionary learning of our proposed CDL algorithm, we sample 480 salient and non-salient image patches of size $80 \times 80$

Table 3: Parameter settings in our method

| Process | Parameter Description | Value |
|---------|---------------------|-------|
| Dictionary learning | Training patch size | $16 \times 16$ |
| | Dictionary atom size $m$ | $256 \times 1$ |
| | Number of atoms in the dictionary $k$ | $1 \times 1024$ |
| | Regularisation parameter $\lambda_1$ | 0.02 |
| | Regularisation parameter $\lambda_2$ | 0.02 |
| Saliency detection | Scalar parameter $\eta_{\mathbf{A}}$ | 1 |
| | Scalar parameter $\eta_{\mathbf{R}}$ | 1 |
| | Positive value $\varphi$ | 0.001 |



(a)                                    (b)

Figure 5: Quantitative comparison of various values of $\lambda_2$. (a) S-measure, adaptive F-measure and E-measure values. (b) MAE values.

from the training set as training patches. During of dictionary training, following the empirical settings in article [17], we down-sample the training patch into an image patch of size $16 \times 16$ as the input, so that the number of pixels $m$ of the learned dictionary atom is 256, and the number of atoms in the dictionary $k$ is set to $4 \times m$ . In Eq. 5, the regularization parameter $\lambda_1$ is set to $1.2/\sqrt{m}$ to weight the reconstruction error and sparsity, and the learning rate $\sigma$ is set to 0.02 in dictionary learning (Eq. 12) to obtain a more discriminative dictionary. Further, we use various values $\lambda 0 = 0.001, \lambda 1 = 0.005, \lambda 2 = 0.01, \lambda 3 = 0.02, \lambda 4 = 0.03, \lambda 5 = 0.04, \lambda 6 = 0.05, \lambda 7 = 0.06, \lambda 8 = 0.07, \lambda 9 = 0.09, \lambda 10 = 0.1$ for testing with respect to $\lambda_2$. The experimental results are shown in Fig. 5, according to which the parameter $\lambda_2$ in Eq. 12 is adjusted to 0.05.
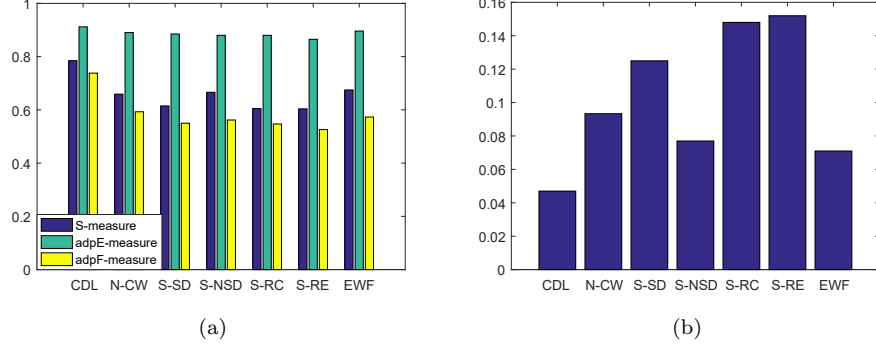
Figure 6: Quantitative comparison of various verification methods. (a) S-measure, adaptive F-measure and E-measure values. (b) MAE values.

## 4.2. Verification and analysis

In this subsection, we first use the 2RSOD dataset to demonstrate the effectiveness of the CDL-based model, then compare our proposed model with state-of-the-art methods on four datasets, and finally analyze some failure cases in our method.

### 4.2.1. Effectiveness analysis based on CDL model

In this subsection, we analyze and verify the effectiveness of the proposed CDL-based saliency detection model on the 2RSOD dataset from the following four aspects:

**A. Effectiveness of contrast-weighted terms**

The contrast-weighted term in Eq. 7 is used to optimize the atoms update during the learning process of the discriminant dictionary. We set it to 1 to verify its effectiveness, as shown in N-CW in Fig. 6 (a).

**B. Effectiveness of constructing a discriminant dictionary**

Discriminant dictionaries have their own features for the SR of images. To verify the effectiveness of the discriminant dictionary for saliency detection, we use the single salient dictionary (S-SD) or single non-salient dictionary (S-NSD) for saliency detection.

**C. The validity of significance measures for joint representation coefficients and reconstruction errors**

To improve the expression of outliers in the coefficients of the SR, we combine representation coefficients and reconstruction errors as a measure of saliency detection. To verify the effectiveness of the joint saliency measure, we use the single representation coefficient (S-RC) and the single reconstruction error (S-RE) as the saliency measurements.

**D. Effectiveness of saliency map fusion method based on global gradient optimization**

To improve the use of correct information in multiple saliency maps, we proposed a saliency map fusion method based on global gradient optimization. We compare the saliency of this proposed optimization method with the equal weight fusion (EWF) method as a verification of the effectiveness of our method.

Fig. 6 shows that the proposed saliency detection method based on CDL is superior to the above effectiveness verification methods, in terms of several evaluation metrics of. The figure also shows the importance and contribution of the various parts that make up the proposed method.

*4.2.2. Comparison with state-of-the-art methods*

We compare the proposed algorithm with nine state-of-the-art SOD methods, including three traditional methods (LPS [51] DSG [52], WMR [30]), three methods related to SR (SMD [25], RSR-LC [27], RDR [26]), and three latest deep learning-based methods (BMPM [34], BASNet [35], PAGE [36]). All results are either generated by the source code or provided by the author.

**A. Visual comparison**

As shown in Fig. 7, most of the comparison methods perform poorly on 2RSOD. On the other hand, the proposed method is competitive on three natural image datasets. Further observation shows that, for the images with simple backgrounds and prominent foregrounds (for example: the second and sixth rows in Fig. 7), all methods have better detection results. However, for images with complex backgrounds (for example: the first row in Fig. 7) that

contain shadows, buildings, and so on, the comparison methods do not have satisfactory results. Moreover, because most of the salient road regions of the images in the 2RSOD dataset are linked to image boundaries, the accuracies of saliency detection methods based on boundary priors (for example, WMR [30], SMD [25], RSR-LC [27], RDR [26]) were also affected. In contrast, the proposed method can effectively separate the salient object from the background and obtains good detection results for images with complex scenes or similar foreground and background.



Figure 7: Visual comparisons of various methods.

### B. Quantitative comparison

To fully compare the proposed method with the above models, the detailed experimental results in terms of four metrics are listed in Table 4. In addition, Fig. 8 shows the standard PR curves and the F-measure curves on the four datasets, which can be used to evaluate the holistic performance of models. As shown in the above experimental results, our proposed method is highly competitive under all six metrics, especially on the 2RSOD dataset. At the same time, the method proposed in this paper is significantly better than the saliency detection method related to SR.

21

Table 4: Quantitative evaluation. The mean F-measure, S-measure and MAE of different saliency detection methods on 2RSOD and three benchmark datasets. The best four results are highlighted in red, blue, green and purple. † & ‡ denote methods based on SR and deep learning. "PAS-S" & "DUT-O" represent datasets PASCAL-S and DUT-OMRON.

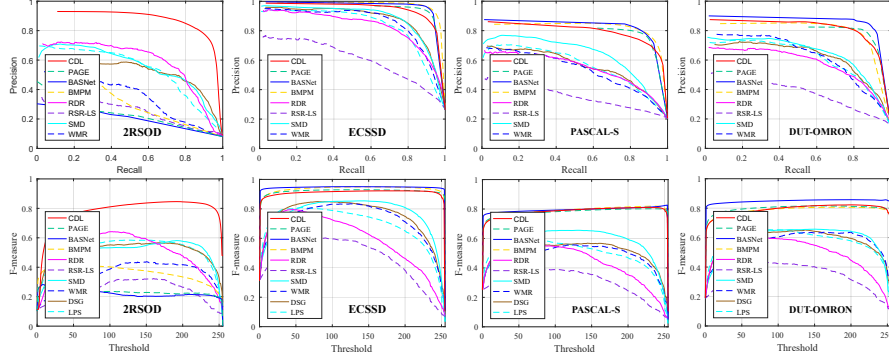| | Metric | LPS [51] | DSG [52] | WMR [30] | SMD† [25] | SRS-LC† [27] | RDR† [26] | BMPM‡ [34] | PAGE‡ [36] | BASNet‡ [36] | CDL (ours) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 2RSOD | $\mathcal{M}\downarrow$ | 0.114 | 0.146 | 0.237 | 0.154 | 0.229 | 0.123 | 0.172 | 0.169 | 0.128 | 0.063 |
| | $F_\beta\uparrow$ | 0.493 | 0.492 | 0.326 | 0.468 | 0.224 | 0.474 | 0.442 | 0.403 | 0.485 | 0.794 |
| | $S_\alpha\uparrow$ | 0.660 | 0.641 | 0.530 | 0.633 | 0.503 | 0.646 | 0.607 | 0.587 | 0.628 | 0.813 |
| | $E_\xi\uparrow$ | 0.811 | 0.709 | 0.629 | 0.772 | 0.729 | 0.818 | 0.734 | 0.701 | 0.743 | 0.901 |
| ECSSD [46] | $\mathcal{M}\downarrow$ | 0.169 | 0.146 | 0.162 | 0.141 | 0.257 | 0.198 | 0.044 | 0.042 | 0.037 | 0.061 |
| | $F_\beta\uparrow$ | 0.629 | 0.701 | 0.669 | 0.725 | 0.422 | 0.549 | 0.894 | 0.906 | 0.880 | 0.879 |
| | $S_\alpha\uparrow$ | 0.700 | 0.773 | 0.754 | 0.795 | 0.592 | 0.641 | 0.911 | 0.912 | 0.916 | 0.884 |
| | $E_\xi\uparrow$ | 0.768 | 0.823 | 0.820 | 0.839 | 0.702 | 0.756 | 0.914 | 0.920 | 0.921 | 0.898 |
| PAS-S [47] | $\mathcal{M}\downarrow$ | 0.203 | 0.231 | 0.249 | 0.198 | 0.276 | 0.218 | 0.037 | 0.077 | 0.076 | 0.088 |
| | $F_\beta\uparrow$ | 0.425 | 0.446 | 0.428 | 0.514 | 0.269 | 0.382 | 0.803 | 0.810 | 0.775 | 0.765 |
| | $S_\alpha\uparrow$ | 0.580 | 0.595 | 0.587 | 0.661 | 0.496 | 0.564 | 0.840 | 0.835 | 0.832 | 0.823 |
| | $E_\xi\uparrow$ | 0.677 | 0.694 | 0.674 | 0.733 | 0.646 | 0.686 | 0.838 | 0.841 | 0.847 | 0.844 |
| DUT-O [48] | $\mathcal{M}\downarrow$ | 0.135 | 0.180 | 0.197 | 0.160 | 0.230 | 0.161 | 0.063 | 0.052 | 0.048 | 0.061 |
| | $F_\beta\uparrow$ | 0.530 | 0.520 | 0.504 | 0.537 | 0.296 | 0.442 | 0.698 | 0.777 | 0.791 | 0.762 |
| | $S_\alpha\uparrow$ | 0.678 | 0.663 | 0.650 | 0.690 | 0.546 | 0.624 | 0.809 | 0.854 | 0.866 | 0.840 |
| | $E_\xi\uparrow$ | 0.748 | 0.760 | 0.729 | 0.746 | 0.659 | 0.731 | 0.839 | 0.869 | 0.884 | 0.860 |



Figure 8: Performance comparison with nine state-of-the-art methods over four datasets. The first row shows a comparison of precision-recall curves. The second row shows a comparison of F-measure curves over different thresholds.

## C. Computational complexity comparison

To demonstrate the computational efficiency of the proposed method, we test the average execution time of several state-of-the-art methods and the proposed method on the 2RSOD dataset. These methods are run on a desktop with an Intel Core i7-7700 CPU and RTX 2070 GPU. As shown in Table 5, the efficiency of the CDL method based on Matlab programming can reach the average of the

Table 5: Average execution time of several methods

| Method | LPS [51] | DSG [52] | WMR [30] | SMD† [25] | SRS-LC† [27] | RDR† [26] | BMPM‡ [34] | PAGE‡ [36] | BASNet‡ [36] | CDL (ours) |
|---|---|---|---|---|---|---|---|---|---|---|
| **Time(s)** | 2.21 | 0.83 | 1.56 | 0.91 | 3.32 | 2.48 | 1.26 | 0.12 | 0.68 | 1.52 |

comparison method.

### 4.2.3. Failure cases

Although the proposed method can accurately detect most salient road regions, there are still some limitations. Fig. 9 shows that when an image contains regions with similar appearance to the road (such as roofs, or farmland), our proposed method incorrectly marks the background regions as the foreground. Also, the places where the road regions would be interrupted are shown in the third column of Fig. 9, which is inconsistent with the fact that the road has connectivity. On the other hand, the saliency of the road should be regional and overall, but as shown in the first and second columns of Fig. 9, there are many scattered points with high saliency values in the inspection results. Through the above analysis, we can construct a more robust dictionary to overcome these problems by combining the semantic information of the image (similar to the work of [53] and [23]) and the feature information of the target, which is one of our future works.

## 5. Conclusion

In this paper, we propose a novel saliency detection method for RS images based on SR. According to the characteristics of salient and non-salient regions, our method uses the proposed online discriminant dictionary learning algorithm to introduce contrast-weighted items into the dictionary learning process to construct a discriminant dictionary based on optimized contrast weighted atoms. Under the discriminant dictionary, we combine the representation coefficients and reconstruction errors of image blocks as saliency detection metrics to generate multiple saliency maps. Considering the complementary informa-
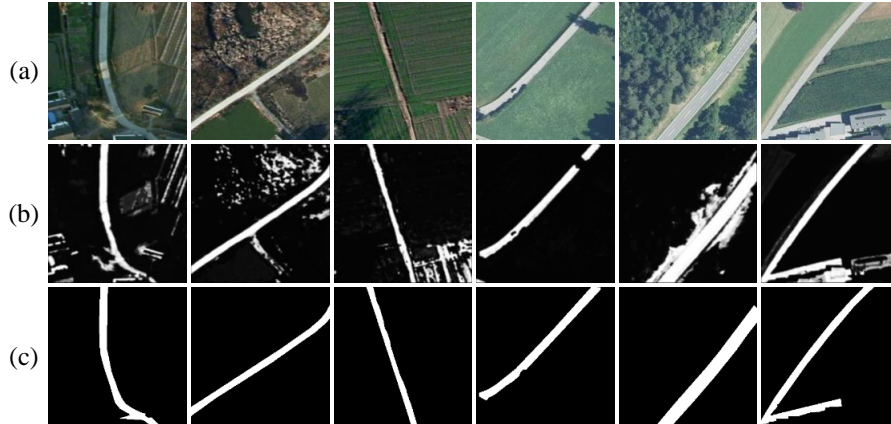
Figure 9: Failure cases of proposed method. (a) Original images. (b) Saliency maps obtained by the proposed method. (c) Ground truth.

.

tion between saliency maps, we propose a saliency map fusion method based on global gradient optimization to integrate multiple saliency maps, which further improves the use of important information from these saliency maps. In addition, we collected and annotated a dataset containing 300 optical RS images. Qualitative, quantitative and ablation experiments on this dataset verify the effectiveness of the proposed method. However, we find that the detection method may fail if an image contains high-contrast areas or has areas similar to the foreground, and the efficiency of the algorithm needs to be improved.

In future work, we will combine the semantic information of the scene and the feature information of the salient object, and develop a more accurate color dictionary to improve the robustness of the multi-class saliency detection. Further, due to the successful use of depth information in SOD, we will explore its application in stereo paired RS data. Inspired by recent work [54] and considering the large-scale and final processing structure of RS images, we also plan to introduce the structure co-occurrence texture (scoot) as a perceptual metric for future SOD work.

## Acknowledgments

## References

[1] A. Borji, M.-M. Cheng, H. Jiang, J. ru Li, Salient object detection: A survey, Computational Visual Media 5 (2019) 117–150.

[2] J. M. Wolfe, T. S. Horowitz, What attributes guide the deployment of visual attention and how do they do it?, Nature reviews neuroscience 5 (6) (2004) 495.

[3] S. K. Mannan, C. Kennard, M. Husain, The role of visual salience in directing eye movements in visual object agnosia, Current biology 19 (6) (2009) R247–R248.

[4] D.-P. Fan, Z. Lin, G.-P. Ji, D. Zhang, H. Fu, M.-M. Cheng, Taking a deeper look at the co-salient object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2020.

[5] A. Borji, D. N. Sihite, L. Itti, Quantitative analysis of human-model agreement in visual saliency modeling: A comparative study, IEEE Transactions on Image Processing 22 (1) (2012) 55–69.

[6] N. Tong, H. Lu, Y. Zhang, X. Ruan, Salient object detection via global and local cues, Pattern Recognition 48 (10) (2015) 3258–3267.

[7] D.-P. Fan, W. Wang, M.-M. Cheng, J. Shen, Shifting more attention to video salient object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 8554–8564.

[8] Q. Fan, C. Qi, Saliency detection based on global and local short-term sparse representation, Neurocomputing 175 (2016) 81–89.

[9] Q. Cai, H. Liu, Y. Qian, S. Zhou, X. Duan, Y.-H. Yang, Saliency-guided level set model for automatic object segmentation, Pattern Recognition 93 (2019) 147–163.

[10] C. F. Flores, A. Gonzalez-Garcia, J. van de Weijer, B. Raducanu, Saliency for fine-grained object recognition in domains with scarce training data, Pattern Recognition 94 (2019) 62–73.

[11] D.-P. Fan, G.-P. Ji, G. Sun, M.-M. Cheng, J. Shen, L. Shao, Camouflaged object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2020.

[12] Y. Zheng, L. Jiao, H. Liu, X. Zhang, B. Hou, S. Wang, Unsupervised saliency-guided sar image change detection, Pattern Recognition 61 (2017) 309–326.

[13] E. Li, S. Xu, W. Meng, X. Zhang, Building extraction from remotely sensed images by integrating saliency cue, IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing 10 (3) (2016) 906–919.

[14] Z. Liu, D. Zhao, Z. Shi, Z. Jiang, Unsupervised saliency model with color markov chain for oil tank detection, Remote Sensing 11 (9) (2019) 1089.

[15] A. Garcia-Diaz, V. Leboran, X. R. Fdez-Vidal, X. M. Pardo, On the relationship between optical variability, visual saliency, and eye fixations: A computational approach, Journal of vision 12 (6) (2012) 17–17.

[16] B. A. Olshausen, D. J. Field, Sparse coding of sensory inputs, Current opinion in neurobiology 14 (4) (2004) 481–487.

[17] J. Mairal, F. Bach, J. Ponce, G. Sapiro, Online dictionary learning for sparse coding, in: Proceedings of the 26th annual international conference on machine learning, ACM, 2009, pp. 689–696.

[18] X. Li, H. Lu, L. Zhang, X. Ruan, M.-H. Yang, Saliency detection via dense and sparse reconstruction, in: Proceedings of the IEEE International Conference on Computer Vision, 2013, pp. 2976–2983.

[19] Y. Zhang, X. Wang, X. Xie, Y. Li, Salient object detection via recursive sparse representation, Remote Sensing 10 (4) (2018) 652.

[20] B. Han, H. Zhu, Y. Ding, Bottom-up saliency based on weighted sparse coding residual, in: Proceedings of the 19th ACM international conference on Multimedia, ACM, 2011, pp. 1117–1120.

[21] J. Yan, M. Zhu, H. Liu, Y. Liu, Visual saliency detection via sparsity pursuit, IEEE Signal Processing Letters 17 (8) (2010) 739–742.

[22] X. Shen, Y. Wu, A unified approach to salient object detection via low rank matrix recovery, in: 2012 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2012, pp. 853–860.

[23] B. Li, C. Yuan, W. Xiong, W. Hu, H. Peng, X. Ding, S. Maybank, Multi-view multi-instance learning based on joint sparse representation and multi-view dictionary learning, IEEE transactions on pattern analysis and machine intelligence 39 (12) (2017) 2554–2560.

[24] Y. Wei, F. Wen, W. Zhu, J. Sun, Geodesic saliency using background priors, in: European conference on computer vision, Springer, 2012, pp. 29–42.

[25] H. Peng, B. Li, H. Ling, W. Hu, W. Xiong, S. J. Maybank, Salient object detection via structured matrix decomposition, IEEE transactions on pattern analysis and machine intelligence 39 (4) (2016) 818–832.

[26] H. Xiao, W. Ren, W. Wang, Y. Liu, M. Zhang, Salient object detection via robust dictionary representation, Multimedia Tools and Applications 77 (3) (2018) 3317–3337.

[27] L. Yi, Z. Qiang, H. Jungong, W. Long, Salient object detection employing robust sparse representation and local consistency, Image and Vision Computing 69 (2018) 155–167.

[28] M. Xu, L. Jiang, Z. Ye, Z. Wang, Bottom-up saliency detection with sparse representation of learnt texture atoms, Pattern Recognition 60 (2016) 348–360.

[29] Q. Zhang, Z. Huo, Y. Liu, Y. Pan, C. Shan, J. Han, Salient object detection employing a local tree-structured low-rank representation and foreground consistency, Pattern Recognition 92 (2019) 119–134.

[30] X. Zhu, C. Tang, P. Wang, H. Xu, M. Wang, J. Chen, J. Tian, Saliency detection via affinity graph learning and weighted manifold ranking, Neurocomputing 312 (2018) 239–250.

[31] D.-P. Fan, M.-M. Cheng, J.-J. Liu, S.-H. Gao, Q. Hou, A. Borji, Salient objects in clutter: Bringing salient object detection to the foreground, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 186–202.

[32] Y. Liu, J. Han, Q. Zhang, L. Wang, Salient object detection via two-stage graphs, IEEE Transactions on Circuits and Systems for Video Technology 29 (4) (2018) 1023–1037.

[33] X. Yao, J. Han, L. Guo, S. Bu, Z. Liu, A coarse-to-fine model for airport detection from remote sensing images using target-oriented visual saliency and crf, Neurocomputing 164 (2015) 162–172.

[34] L. Zhang, J. Dai, H. Lu, Y. He, G. Wang, A bi-directional message passing model for salient object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 1741–1750.

[35] X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, M. Jagersand, Basnet: Boundary-aware salient object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 7479–7489.

[36] W. Wang, S. Zhao, J. Shen, S. C. Hoi, A. Borji, Salient object detection with pyramid attention and salient edges, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 1448–1457.

[37] K. Fu, D.-P. Fan, G.-P. Ji, Q. Zhao, Jl-dcf: Joint learning and densely-cooperative fusion framework for rgb-d salient object detection, in: Pro-

ceedings of the IEEE conference on computer vision and pattern recognition, 2020.

[38] J. Zhang, D.-P. Fan, Y. Dai, S. Anwar, F. Sadat Saleh, T. Zhang, N. Barnes, Uc-net: Uncertainty inspired rgb-d saliency detection via conditional variational autoencoders, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2020.

[39] D.-P. Fan, Z. Lin, J.-X. Zhao, Y. Liu, Z. Zhang, Q. Hou, M. Zhu, M.-M. Cheng, Rethinking rgb-d salient object detection: Models, datasets, and large-scale benchmarks, IEEE Transactions on Neural Networks and Learning Systems (2020).

[40] J.-X. Zhao, Y. Cao, D.-P. Fan, M.-M. Cheng, X.-Y. Li, L. Zhang, Contrast prior and fluid pyramid integration for rgbd salient object detection, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[41] M.-M. Cheng, N. J. Mitra, X. Huang, P. H. Torr, S.-M. Hu, Global contrast based salient region detection, IEEE Transactions on Pattern Analysis and Machine Intelligence 37 (3) (2014) 569–582.

[42] M. R. Osborne, B. Presnell, B. A. Turlach, A new approach to variable selection in least squares problems, IMA journal of numerical analysis 20 (3) (2000) 389–403.

[43] S. J. Wright, Coordinate descent algorithms, Mathematical Programming 151 (1) (2015) 3–34.

[44] G.-S. Xia, X. Bai, J. Ding, Z. Zhu, S. Belongie, J. Luo, M. Datcu, M. Pelillo, L. Zhang, Dota: A large-scale dataset for object detection in aerial images, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 3974–3983.

[45] G. Cheng, P. Zhou, J. Han, Learning rotation-invariant convolutional neural networks for object detection in vhr optical remote sensing images, IEEE Transactions on Geoscience and Remote Sensing 54 (12) (2016) 7405–7415.

[46] Q. Yan, L. Xu, J. Shi, J. Jia, Hierarchical saliency detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2013, pp. 1155–1162.

[47] Y. Li, X. Hou, C. Koch, J. M. Rehg, A. L. Yuille, The secrets of salient object segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 280–287.

[48] C. Yang, L. Zhang, H. Lu, X. Ruan, M.-H. Yang, Saliency detection via graph-based manifold ranking, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2013, pp. 3166–3173.

[49] D.-P. Fan, M.-M. Cheng, Y. Liu, T. Li, A. Borji, Structure-measure: A new way to evaluate foreground maps, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 4548–4557.

[50] D.-P. Fan, C. Gong, Y. Cao, B. Ren, M.-M. Cheng, A. Borji, Enhanced-alignment Measure for Binary Foreground Map Evaluation, in: International Joint Conference on Artificial Intelligence (IJCAI), 2018, pp. 698–704, http://dpfan.net/e-measure/.

[51] H. Li, H. Lu, Z. Lin, X. Shen, B. Price, Inner and inter label propagation: salient object detection in the wild, IEEE Transactions on Image Processing 24 (10) (2015) 3176–3186.

[52] L. Zhou, Z. Yang, Z. Zhou, D. Hu, Salient region detection using diffusion process on a two-layer sparse graph, IEEE Transactions on Image Processing 26 (12) (2017) 5882–5894.

[53] Y. Liu, M.-M. Cheng, D.-P. Fan, L. Zhang, J. Bian, D. Tao, Semantic edge detection with diverse deep supervision, arXiv preprint arXiv:1804.02864 (2018).

[54] D.-P. Fan, S. Zhang, Y.-H. Wu, Y. Liu, M.-M. Cheng, B. Ren, P. L. Rosin, R. Ji, Scoot: A perceptual metric for facial sketches, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 5612–5622.