

SpineOne: A One-Stage Detection Framework for Degenerative Discs and Vertebrae

Jiabo He^{1,2}, Wei Liu², Yu Wang², Xingjun Ma³, Xian-Sheng Hua²

¹School of Computing and Information Systems, The University of Melbourne, Australia

²DAMO Academy, Alibaba Group, China

³School of Information Technology, Deakin University, Australia

jiaboh@student.unimelb.edu.au, daniel.ma@deakin.edu.au, {vivi.lw, tonggou.wangyu, xiansheng.hxs}@alibaba-inc.com

arXiv:2110.15082v2 [cs.CV] 10 Nov 2021

Abstract—Spinal degeneration plagues many elders, office workers, and even the younger generations. Effective pharomic or surgical interventions can help relieve degenerative spine conditions. However, the traditional diagnosis procedure is often too laborious. Clinical experts need to detect discs and vertebrae from spinal magnetic resonance imaging (MRI) or computed tomography (CT) images as a preliminary step to perform pathological diagnosis or preoperative evaluation. Machine learning systems have been developed to aid this procedure generally following a two-stage methodology: first perform anatomical localization, then pathological classification. Towards more efficient and accurate diagnosis, we propose a one-stage detection framework termed SpineOne to simultaneously localize and classify degenerative discs and vertebrae from MRI slices. SpineOne is built upon the following three key techniques: 1) a new design of the keypoint heatmap to facilitate simultaneous keypoint localization and classification; 2) the use of attention modules to better differentiate the representations between discs and vertebrae; and 3) a novel gradient-guided objective association mechanism to associate multiple learning objectives at the later training stage. Empirical results on the Spinal Disease Intelligent Diagnosis Tianchi Competition (SDID-TC) dataset of 550 exams demonstrate that our approach surpasses existing methods by a large margin.

Index Terms—Magnetic resonance imaging, one-stage detection, discs and vertebrae, spinal degeneration.

I. INTRODUCTION

SPINAL diseases have become increasingly common nowadays, among which the degeneration of spines is nearly inevitable with aging. Degenerative spine conditions involve the gradual loss of normal structures and functions over time, which may be caused by aging, tumors, infections and arthritis [1]. Magnetic resonance imaging (MRI) and computed tomography (CT) techniques are used to visualize the anatomical structures of the spine before pathological diagnosis and treatment. The lumbar spine consists of 5 discs and 5 vertebrae (Fig. 1). Discs act as shock absorbers between vertebrae and both of them can involve degenerative conditions [2].

Degenerative spine conditions can be relieved by pharomic or surgical interventions. Such labor-intensive interventions rely on clinical experts to manually identify degenerative discs and vertebrae from spinal MRI or CT images. Traditional image processing systems have been developed to help localize discs

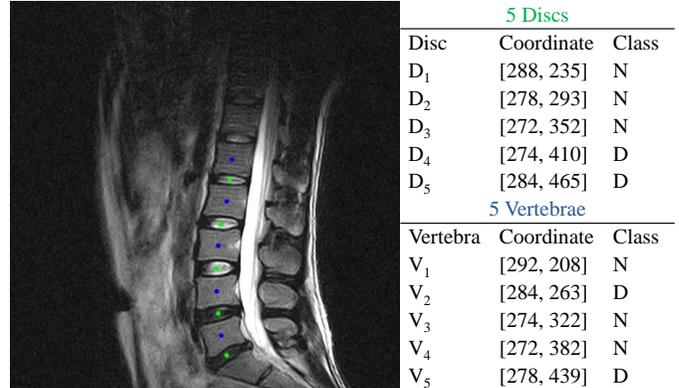


Fig. 1. Illustration of 5 discs (centroids in green) and 5 vertebrae (centroids in blue) from the MRI slice of the lumbar spine. D_1 , D_2 , D_3 , D_4 , and D_5 denote the disc L_1-L_2 , L_2-L_3 , L_3-L_4 , L_4-L_5 , and L_5-S_1 , respectively; V_1 , V_2 , V_3 , V_4 , and V_5 denote the vertebra L_1 , L_2 , L_3 , L_4 , and L_5 , respectively; where L_1-5 are the five lumbar levels from top to bottom and S_1 is the sacral spine. N and D denote Normal and Degenerative classes.

and vertebrae via an automated process [3]–[5]. Convolutional neural networks (CNNs) have also been introduced to detect and segment discs and vertebrae from spinal MRIs [6]–[8]. For example, SpineNet generates slice-level (i.e., image-level) radiological scores for discs [6] while Spine-GAN segments discs, vertebrae, and neural foramen with pixel-level results [7]. Different from above existing works, we localize discs and vertebrae by predicting the locations of their centroids (i.e., keypoints), as well as generating their anatomical-structure-level (i.e., keypoint-level) degenerative categories. We will further show that our task enables minimal annotation effort from experts compared with existing segmentation tasks.

There are three main challenges of applying existing methods to our task: 1) there are multiple centroids (5 discs and 5 vertebrae) and multiple pathological classes (2 classes for either discs or vertebrae); 2) discs and vertebrae need to be diagnosed simultaneously, despite their high inter-structural similarity and intra-structural variation; and 3) there are multiple learning objectives including localization and classification (for each disc and vertebra) and the localization involves both keypoint localization and offset regression. To address these challenges, we propose *SpineOne*, an efficient one-stage framework that can automatically localize and classify

This research was done when Jiabo He interned with the DAMO Academy, Alibaba Group. A short version of this paper has been accepted by IEEE BIBM 2021. Correspondence to Xingjun Ma and Yu Wang.

degenerative discs and vertebrae from MRIs. SpineOne is more friendly to clinical experts than two-stage frameworks, since it can provide the localization and classification results at the same time. SpineOne takes MRI slices of a lumbar spine as inputs and outputs centroids and offsets for detected discs and vertebrae, and at the same time, outputs the pixel-wise probability of each disc or vertebra being degenerative. Therefore, SpineOne is able to help clinical experts relieve from large amounts of the workload by serving as an assistant, making efficient and accurate diagnosis.

Specifically, based on CNNs, SpineOne addresses the above three challenges using three novel techniques. **First**, a one-channel-per-class (OCPC) keypoint heatmap is designed to promote simultaneous localization and classification without introducing additional heads or output channels into the network. Our OCPC design is advantageous as the centroids of discs and vertebrae are spatially separated following physiological rigidity, which can help capture the geometrical and classification correlations among keypoints (Fig. 1). **Second**, we introduce dual self-attention modules to adaptively integrate similar features at different scales – an inspiration from human experts who would inspect MRI slices from a global view with visual acuity. Specifically, we use 1) a position attention module (PAM) to aggregate the feature at each position by a weighted sum of features over all positions of the image, and 2) a channel attention module (CAM) to integrate features across channels. **Third**, we introduce a novel gradient-guided objective association (OA) mechanism to adaptively associate two types of objectives, i.e., using the gradient of the loss w.r.t. the heatmap to guide the learning of the offset. OA is a generic technique to boost the learning of the main objectives in multi-head models at the later training stage.

In summary, our main contributions are:

- We propose a novel one-channel-per-class (OCPC) keypoint heatmap design to facilitate simultaneous keypoint localization and classification. Our OCPC heatmap compacts all keypoints of the same class into the same channel given they are spatially separated.
- We introduce dual self-attention modules for the learning of more distinguishable representations between discs and vertebrae: a position attention module (PAM) to model inner-image spatial interaction, and a channel attention module (CAM) to model inter-channel correlation.
- We propose a novel gradient-guided objective association (OA) mechanism to explicitly connect heatmap learning with offset learning. This makes the learning of the heatmap more effective, which is the main objective of keypoint localization and classification.
- We integrate above three techniques into one novel one-stage detection framework SpineOne, and empirically show, on the Spinal Disease Intelligent Diagnosis Tianchi Competition dataset [9], that SpineOne can outperform the current state-of-the-art methods by a large margin.

II. RELATED WORK

In this section, we first review existing deep learning models developed for spine-related tasks. We then discuss various state-of-the-art methods for keypoint detection.

A. Segmentation, Detection and Pathological Classification on Spines

A number of deep learning models have been proposed for spine-related tasks including anatomical segmentation [10]–[12], detection [8], [13]–[15], and pathological classification [6], [16] from either MRIs or CT images. These works have provided great assistance to clinical experts in osteoporosis assessment, fractures detection and aging process analysis. There have also been several multi-task learning studies on spines. For example, two-stage methods have been proposed to train one network for vertebral segmentation and disc image extraction, and the other network for stenosis grading [17]. Different from these two-stage methods, one-stage methods are able to segment discs, vertebrae, and neural foramen simultaneously using GAN-based models [7]. The sequential conditional reinforcement learning network (SCRL) can also tackle the simultaneous vertebral detection and segmentation from MRIs [18]. 2D CNNs are frequently used to process MRI slices while 3D CNNs are popular for 3D CT images.

While there exists a body of work on spine-related tasks, our work is notably different. Specifically, our task requires to localize 5 discs and 5 vertebrae on the lumbar spine at the same time by localizing their centroids, which minimizes the annotation effort of experts. Furthermore, we classify the degenerative conditions for each disc and vertebra. To the best of our knowledge, none of the existing works has addressed the same task before. In this work, we address this gap by proposing an efficient and accurate one-stage framework.

B. Keypoint Detection

Keypoint detection is fundamental to numerous vision tasks [19]–[21], amongst which human pose estimation and keypoint-based object detection are the two most related topics to our task.

1) *Human pose estimation*: Single-person pose estimation localizes human anatomical keypoints/parts by regressing either spatial joint coordinates [22] or location heatmaps [23], [24]. Multi-person pose estimation can be achieved via either bottom-up [25], [26] or top-down [27], [28] approaches with fully CNNs. Besides fully CNNs, specialized CNN architectures were also proposed for keypoint detection. For example, the stacked hourglass network repeats the bottom-up and top-down processing in conjunction with intermediate supervision [29]. The HRNet maintains high-resolution representations through the entire training process using high-to-low resolution sub-networks [30]. Keypoint detection requires rich spatial correlations and contextual information captured from both high- and low-level representations [31], [32]. As such, keypoint detection can potentially boost other tasks via single-shot approaches, e.g., instance segmentation, semantic

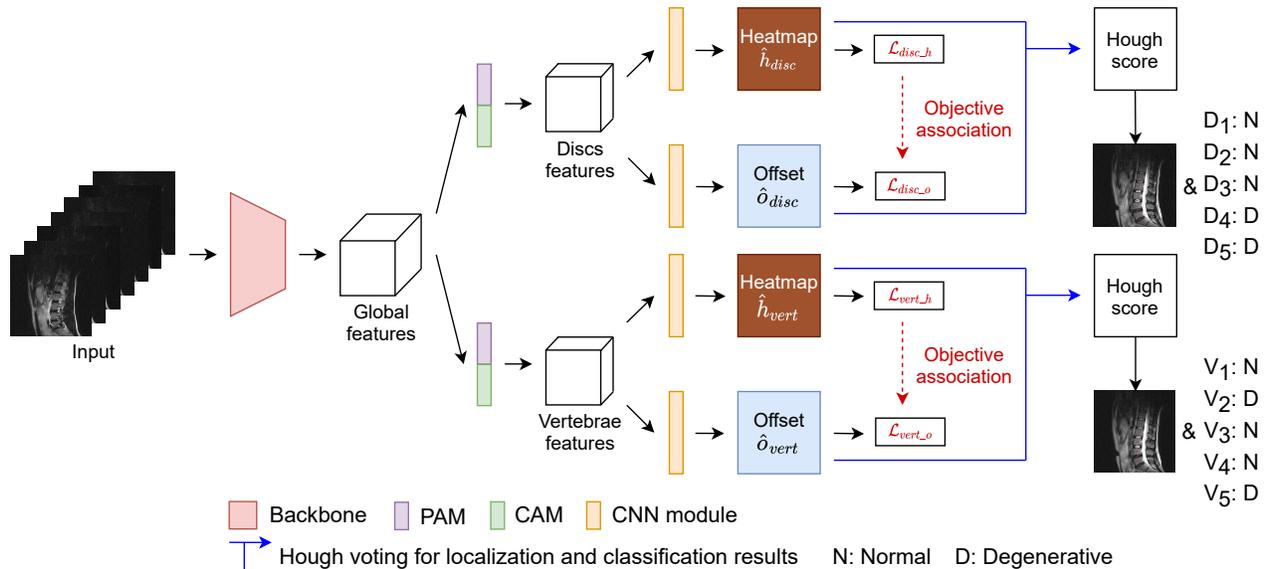


Fig. 2. Overview of our proposed framework SpineOne. One example output of our framework can be found in Fig. 1. The structures of PAM and CAM modules are shown in Fig. 3.

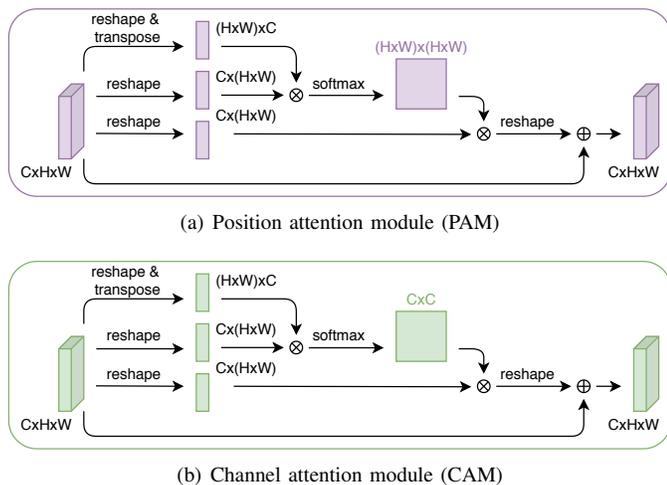


Fig. 3. Structure details of PAM and CAM. \otimes : matrix multiplication; \oplus : element-wise addition.

segmentation, and even image parsing (i.e., panoptic segmentation) [33]–[35]. In this work, state-of-the-art architectures for keypoint detection are leveraged at the first stage of two-stage baselines, followed by classifiers for those positive/detected keypoints at the second stage.

2) *Keypoint-based object detection*: These methods detect objects based on the keypoint heatmaps by one-stage, anchor-free detectors. CornetNet [36], CenterNet1 [37], ExtremeNet [38], and CentripetalNet [39] predict locations of keypoints (corners, centroids, or extreme points) and group them into same bounding boxes if they are geometrically aligned. Compared with these works, we only need to localize the centroids of discs and vertebrae and no grouping is needed. Several works address the simultaneous localization and classification

problem by generating pixel-wise results. CenterNet2 estimates pixel-level categories of objects along with their sizes and offsets [40]. FCOS generates pixel-wise classification, center-ness, and bounding box (top, down, left, right) results using multi-head CNNs [41]. In our task, it is also required to simultaneously localize and classify degenerative conditions of each disc and vertebra. We tackle this challenge by introducing a novel keypoint heatmap to these architectures.

III. PROPOSED ONE-STAGE APPROACH

Problem definition. Let $I \in \mathbb{R}^{n \times H \times W}$ be the sagittal T2 sequence of an input exam with n slices, height H , and width W . SpineOne outputs the keypoint heatmaps (\hat{h}_{disc} and \hat{h}_{vert}) and the corresponding offset maps (\hat{o}_{disc} and \hat{o}_{vert}) for both discs and vertebrae, respectively. The goal is to design and train SpineOne to output consistent \hat{h}_{disc} , \hat{h}_{vert} , \hat{o}_{disc} , and \hat{o}_{vert} with the ground truth h_{disc} , h_{vert} , o_{disc} , and o_{vert} . The ground truth can be generated from the annotations (locations and classes). After postprocessing above outputs, we compute the *precision*, *recall*, and *F1 score* to evaluate the final performance of SpineOne comprehensively.

A. Framework Overview

An ordinary detector has at least three components, i.e., input, backbone, and head [42]. Recent detectors also insert a neck component between the backbone and the head. For clarity, we take the backbone along with the neck (if there is one) as a whole module for feature extraction, fusion, and learning. The proposed SpineOne framework is illustrated in Fig. 2. For either the discs (top) or vertebrae (bottom) branch, traditional two-stage frameworks use sequential output heads with one for localization and the other for classification. SpineOne simplifies this process by using one single head (the ‘Heatmap’ output in Fig. 2) for two purposes, i.e., keypoint

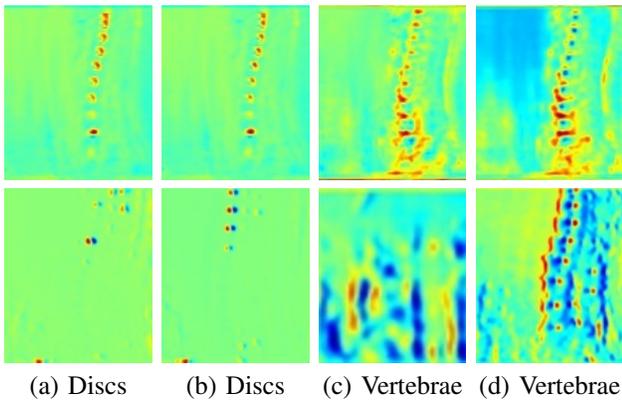


Fig. 4. Visualization of feature maps learned by our framework w/ (top) or w/o (bottom) attention modules. Those features are randomly selected channels of discs features and vertebrae features.

localization and classification. Note that it still needs one additional head for the offset learning, respectively.

We select $n = 7$ slices in the middle of the sagittal T2 sequence of MRIs as inputs, which are the main focus of clinical experts when diagnosing degenerative spines [6]. A backbone is used to extract the global feature map U_{global} , which is then forked by two dual self-attention modules (i.e., PAM and CAM) to learn discs and vertebrae feature representations (U_{disc} and U_{vert}). The two feature representations are then separately fed into CNN modules to generate keypoint heatmaps (\hat{h}_{disc} , and \hat{h}_{vert}), and corresponding offset maps (\hat{o}_{disc} and \hat{o}_{vert}), which is done for discs and vertebrae in parallel. Details of the framework structure is accessible to readers in Table VI. Loss functions are defined between outputs (i.e., \hat{h}_{disc} , \hat{h}_{vert} , \hat{o}_{disc} , and \hat{o}_{vert}) and their ground truth (i.e., h_{disc} , h_{vert} , o_{disc} , and o_{vert}) [43], [44]. Note that we leverage the gradient of the heatmap loss w.r.t. the heatmap to guide the learning of the offset by the objective association (OA) mechanism. The entire network is trained in an end-to-end and supervised fashion. There is also a postprocessing step to produce localization and classification results with Hough voting.

B. Attention Modules

We introduce attention modules for capturing the long-range dependencies among input slices [45]–[49]. Specifically, we use the dual self-attention modules [50] to learn more discriminative disc and vertebra representations. The dual self-attention modules consist of a position attention module (PAM) and a channel attention module (CAM). PAM selectively aggregates the feature at each position by a weighted sum of features at all positions within the image (Fig. 3a), while CAM can notice interdependent channel maps by integrating associated features among all channels (Fig. 3b). Note that the attention maps of PAM and CAM are different in size as feature matrixes are multiplied in reversed orders. It is worth mentioning that PAM is memory hungry as there can be \sim billions of elements in the intermediate attention map of size $(H \times W) \times (H \times W)$ (e.g., there are approximately 0.3

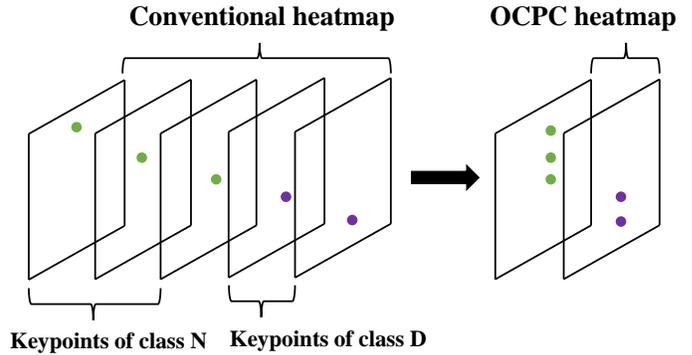


Fig. 5. Comparison between the conventional heatmap and our proposed one-channel-per-class (OCPC) heatmap. Different colors of keypoints represent different classes.

billion elements when $H = W = 128$). Fig. 4 provides a visual comparison between feature maps learned w/ and w/o dual attention modules from MRIs. With attention modules (the top row), the network is able to accurately locate the discs and vertebrae, and pays more attention on the lumbar spine. Moreover, the attention module for discs mainly focuses on regions within discs whereas the module for vertebrae focuses on both regions within vertebrae and their surrounding regions. However, without attention modules, the network may fail to locate discs nor vertebrae, and the attention may be shifted towards unimportant regions. We will also show in our experiments that the final performance can drop considerably without attention modules.

C. OCPC Keypoint Heatmap

The conventional keypoint heatmap has only one positive (i.e., 1) value indicating the location of the ground truth keypoint, while all other locations are zeros. For example, if the heatmap is in size of 640×640 , there will be 409,599 zeros and only 1 value of one. The foreground/background ratio is then $1/409599$ for one keypoint. In order to alleviate the dominance of the background, conventional approaches generate either Gaussian heatmaps [51], [52] or binary heatmaps [33], [53], which involve the processing of local regions (within a radius R) around keypoints and the extension to multiple channels with one channel for each keypoint of each class. This basically means that, for an object with K keypoints and C classes, a heatmap with $K + C$ channels will be generated. These methods force the network to have either redundant output heads or heads with redundant channels, resulting in a complex network architecture that is hard to train. As such, these designs are not suitable for medical applications where a general lack of data hinders the training of complex models. Alternatively, we propose to convert the conventional heatmap to a one-channel-per-class (OCPC) keypoint heatmap (Fig. 5), which not only compacts channels of heads within $C = 2$ channels (the normal class versus the degenerative class), but also captures the geometrical and classification correlations among keypoints.

Let \mathbf{x} be the middle slice of the sagittal T2 sequence, with x_i ($i = 1, \dots, N$) being the positional index of the slice. N is the total number of pixels. Let y_k be the k^{th} keypoint, and $\mathcal{D}_R(y_k) = \{x_i : \|x_i - y_k\| \leq R\}$ be a disk of radius R centered of y_k . For each keypoint y_k of class c , its OCPC heatmap is defined as follows:

$$p_{k,c}(x_i) = \begin{cases} 1, & \text{if } x_i \in \mathcal{D}_R(y_k), \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

This will produce a dense binary OCPC heatmap with $C = 2$ channels for all keypoints with integrated class information. We set R to 6 pixels (i.e., $0.4375 \times 6 = 2.625\text{mm}$ after spacing alignment and resizing all slices to 640×640 pixels) such that disks do not overlap with each other. The foreground/background ratio thus increases from $1/409599$ to approximately $1/4403$ for one keypoint, which can further increase with multiple keypoints of the same class generated in the same channel. OCPC is a compact and generic heatmap design that can be readily applied to any applications where the keypoints are spatially separated from each other.

D. Offset and Hough Voting

1) *Short-range offset*: There are three types of offsets to be considered: 1) short-range, 2) mid-range, and 3) long-range offsets [33]. In our framework, we adopt the short-range offset to improve the keypoint localization. The mid-range offset associates keypoints of the same instance while the long-range offset calculates distances between the keypoint and all points within its instance for other tasks, e.g., instance segmentation. We thus do not consider the latter two offsets in our task. The short-range offset with $2C = 4$ channels (2 channels of coordinate x and y for each keypoint of each class) that points from an image position x_i to the k^{th} keypoint of class c is defined as:

$$S_{k,c}(x_i) = y_k - x_i, \quad x_i \in \mathcal{D}_R(y_k). \quad (2)$$

2) *Hough voting*: We further aggregate the keypoint heatmap and the short-range offset map into a Hough score map $h_{k,c}(x)$ for each class c [33]:

$$h_{k,c}(x) = \frac{1}{\pi R^2} \sum_{i=1}^N p_{k,c}(x_i) B(x_i + S_{k,c}(x_i) - x), \quad (3)$$

where $p_{k,c}(x_i)$ is the sigmoid probability in channel c in the heatmap and $B(\cdot)$ is the bilinear interpolation kernel for size match between outputs and the ground truth. We sort and select the top five points with the top five highest scores in the Hough score map as the prediction results for discs and vertebrae, respectively. The Hough score map facilitates accurate keypoint localization by considering both keypoint locations and their surrounding pixels. We refer readers to [33] for more details of the Hough voting.

E. Gradient-Guided Objective Association

Our SpineOne framework has 4 output heads with two heads for discs and the other two for vertebrae. We combine the loss functions defined for each head into one total loss:

$$\mathcal{L} = \mathcal{L}_{disc_h} + \mathcal{L}_{disc_o} + \mathcal{L}_{vert_h} + \mathcal{L}_{vert_o}, \quad (4)$$

where the first two terms are for discs (i.e., \mathcal{L}_{disc_h} for keypoint heatmap and \mathcal{L}_{disc_o} for offset) and the last two terms are for vertebrae (i.e., \mathcal{L}_{vert_h} and \mathcal{L}_{vert_o}). We test different sets of hyper-parameters and find no noticeable improvement based on the above loss. The four loss terms are thus naturally balanced in our task and introducing hyper-parameters into the above loss is not necessary. Regarding the exact forms of the loss functions, we use the focal loss ($FL(\hat{p}) = -(1 - \hat{p})^\gamma \log(\hat{p})$, $\gamma > 0$) [54] for heatmap losses, and the L_1 loss ($L_1(\hat{r}) = |\hat{r} - r|$) for offset losses.

One remaining problem is that the heatmap generation is in fact more important than the offset regression because the heatmap provides both localization and classification results while the offset is useful only to the localization task. As such, the learning of the heatmap should lead the learning of the offset. To address this problem, we propose a gradient-guided *objective association* (OA) mechanism that associates heatmap loss terms with offset loss terms. Taking discs for an example, its two loss terms are:

$$\begin{aligned} \mathcal{L}_{disc} &= \mathcal{L}_{disc_h} + \mathcal{L}_{disc_o} \\ &= \mathcal{L}_{disc_h}(\hat{\mathbf{h}}_{disc}, \mathbf{h}_{disc}) + \mathcal{L}_{disc_o}(\hat{\mathbf{o}}_{disc}, \mathbf{o}_{disc}), \end{aligned} \quad (5)$$

where $\hat{\mathbf{h}}_{disc}$ and $\mathbf{h}_{disc} \in \mathbb{R}^{C \times H \times W}$ are predicted and ground truth keypoint heatmaps respectively, while $\hat{\mathbf{o}}_{disc}$ and $\mathbf{o}_{disc} \in \mathbb{R}^{2C \times H \times W}$ are predicted and ground truth short-range offset maps respectively. C is the total number of classes. The proposed OA mechanism associates the above two disc loss terms as follows:

$$\begin{aligned} \mathcal{L}_{disc} &= \mathcal{L}_{disc_h}(\hat{\mathbf{h}}_{disc}, \mathbf{h}_{disc}) + \\ &\quad \mathcal{L}_{disc_o}(\mathbf{R}(\mathbf{1} \oplus \frac{\partial \mathcal{L}_{disc_h}}{\partial \hat{\mathbf{h}}_{disc}}) \odot \hat{\mathbf{o}}_{disc}, \mathbf{o}_{disc}), \end{aligned} \quad (6)$$

where the constant tensor $\mathbf{1} \in \mathbb{R}^{C \times H \times W}$ and the gradient of the heatmap loss w.r.t. the heatmap ($\partial \mathcal{L}_{disc_h} / \partial \hat{\mathbf{h}}_{disc}$) are element-wise added up as a type of weight to differentiate the varying importance of pixels in the heatmap. $\mathbf{R}(\cdot)$ is a kernel that reshapes the size of the input from $\mathbb{R}^{C \times H \times W}$ to $\mathbb{R}^{2C \times H \times W}$ by duplicating channels. \odot denotes element-wise product. We clip each gradient value in $\partial \mathcal{L}_{disc_h} / \partial \hat{\mathbf{h}}_{disc}$ to $[-100, 100]$, then normalize it to $[0, 1]$ for each channel. Similarly, we also associate the two loss terms for vertebrae as follows:

$$\begin{aligned} \mathcal{L}_{vert} &= \mathcal{L}_{vert_h}(\hat{\mathbf{h}}_{vert}, \mathbf{h}_{vert}) + \\ &\quad \mathcal{L}_{vert_o}(\mathbf{R}(\mathbf{1} \oplus \frac{\partial \mathcal{L}_{vert_h}}{\partial \hat{\mathbf{h}}_{vert}}) \odot \hat{\mathbf{o}}_{vert}, \mathbf{o}_{vert}). \end{aligned} \quad (7)$$

The gradient of the heatmap loss w.r.t. the heatmap highlights the importance of each pixel in the keypoint heatmap to both localization and classification tasks. Therefore, the OA mechanism encourages the network to focus more on

TABLE I

EXPERIMENTAL RESULTS OF OUR FRAMEWORK COMPARED WITH STATE-OF-THE-ART FRAMEWORKS ON THE SDID-TC DATASET. ALL EVALUATION RESULTS ARE THE MACRO MEAN OF NORMAL AND DEGENERATIVE CLASSES. THE TOP TWO BEST RESULTS IN EVERY COLUMN ARE **BOLDFACED**.

Framework	Backbone	Discs			Vertebrae		
		Recall	Precision	F1	Recall	Precision	F1
Two-stage	Hourglass-52+Res18	0.800	0.831	0.815	0.773	0.797	0.785
	Hourglass-104+Res18	0.808	0.837	0.822	0.782	0.801	0.791
	HRNet-W32+Res18	0.809	0.839	0.823	0.771	0.790	0.780
	HRNet-W48+Res18	0.816	0.844	0.829	0.788	0.807	0.797
One-stage	DenseNet121 w/ FPN	0.802	0.836	0.818	0.755	0.778	0.766
	DenseNet169 w/ FPN	0.811	0.845	0.827	0.767	0.795	0.781
	FCOS(Res101)	0.814	0.848	0.830	0.770	0.798	0.783
	DeepLabv3+(Res50)	0.813	0.848	0.830	0.763	0.784	0.774
	DeepLabv3+(Res101)	0.819	0.855	0.836	0.774	0.795	0.785
	CenterNet(Res101)	0.817	0.851	0.833	0.774	0.802	0.788
SpineOne	DenseNet121 w/ FPN	0.837	0.869	0.852	0.806	0.827	0.816
	DenseNet169 w/ FPN	0.843	0.875	0.858	0.816	0.837	0.827
	FCOS(Res101)	0.849	0.878	0.863	0.822	0.842	0.832
	DeepLabv3+(Res50)	0.847	0.873	0.860	0.811	0.834	0.823
	DeepLabv3+(Res101)	0.859	0.891	0.874	0.840	0.859	0.849
	CenterNet(Res101)	0.857	0.888	0.872	0.839	0.860	0.849

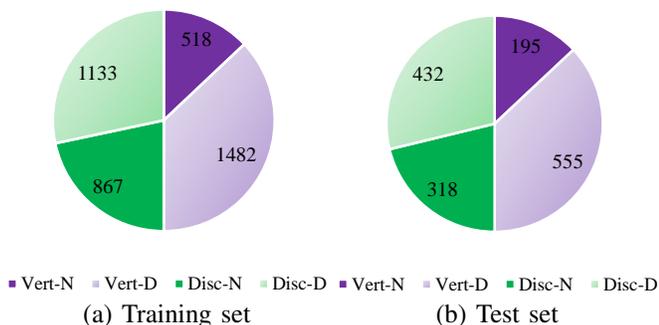


Fig. 6. Statistics of the SDID-TC dataset. Vert-N, Vert-D, Disc-N, and Disc-D denote the numbers of vertebrae being normal or degenerative, discs being normal or degenerative in the dataset, respectively.

important regions learned from keypoint heatmaps, which in turn provides certain guidance for the learning of offsets. OA is hyper-parameter free, thus can be easily incorporated into any existing multi-task learning frameworks to boost the performance of main tasks. We recommend using OA after 75% training epochs when models start to converge.

IV. EXPERIMENTS

A. SDID-TC Dataset

We evaluate our SpineOne framework on the Spinal Disease Intelligent Diagnosis Tianchi Competition (SDID-TC) dataset [9], which inspired researchers to develop efficient and accurate AI frameworks for the purpose of solving medical problems. We randomly split the 550 exams into 400 for training and 150 for test. The distribution of every class is shown in Fig. 6, with mild class imbalance in the dataset. There are both T1 and T2 weighted sequences of MRIs in each exam. We only consider the sagittal T2 sequence as it high-

lights the cerebrospinal fluid (CSF) beside discs and vertebrae, which can help experts identify degenerative conditions more easily based on the brightly-shown deformation of discs and vertebrae, as well as their relative positions to spinal canals. To minimize the annotation work of experts, they were only required to annotate the centroids of discs and vertebrae on the middle slice of sagittal T2, as well as their corresponding degenerative conditions. This style of annotations saved a lot of time compared with annotations of bounding boxes or segmentation masks.

B. Evaluation Metrics

To evaluate our overall diagnosis performance (localization and classification), we first adopt the standard Percentage of Correct Keypoints (PCK) metric in keypoint localization tasks. PCK regards all generated keypoints that fall within a certain threshold distance to the ground truth as positive/detected ones. In our experiments, PCK under $6mm$ (PCK-6) is suggested by clinical experts as the threshold for positive discs and vertebrae. Among all positive centroids, correctly classified (i.e., correct pathological classification) ones are true positive (TP) while wrongly classified ones are false positive (FP). Those missed centroids are false negative (FN). After defining TP, FP, and FN, we further use them to compute the *precision*, *recall*, and *F1 score* as the final evaluation for all frameworks comprehensively.

C. Implementation Details

A number of state-of-the-art methods are compared in our experiments. DeepLab series perform excellently in semantic segmentation, object detection, and panoptic segmentation [34]. Here we use DeepLabv3+ with ResNet50/101 [55] as the backbone, followed by the atrous spatial pyramid pooling

(ASPP) consisting of a 1×1 convolution and three 3×3 convolutions with *rates* = (6, 12, 18). DenseNet121/169 [56] is also compared as a backbone for global feature extraction, followed by feature pyramid networks (FPN) [57] as its neck. In addition, we also compare CenterNet (with ResNet101 as the backbone) [40] and FCOS (with ResNet101 as the backbone) [41] with their heads modified for our task. Different from above adapted one-stage methods, two-stage methods localize keypoints and then classify cropped regions around them subsequently. In order to demonstrate the superiority of SpineOne over two-stage methods, we utilize Hourglass [29] and HRNet [30] for keypoint localization, followed by the ResNet18 as a classifier for those region proposals. Hourglass-52/104 denotes 1 or 2 hourglass modules. HRNet-W32/48 denotes different widths of the high-resolution subnetworks. We also tried deeper ResNets as classifiers but we did not obtain notable improvement.

We set dimensions as $H = W = 128$, $C = 64$ (Fig. 3) for attention modules throughout all models. We set the batch size as 64/16 for models without/with attention modules due to GPU memory constraints. The initial learning rate is set to 0.04/0.01 for different batch sizes and gradually decayed following $lr_i = (1 - i/T)^{0.9}$, at the i^{th} epoch with $T = 300$ epochs in total. We use the Adam [58] optimizer with default settings ($\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$). We select $\gamma = 2$ in the focal loss [54]. In the SDID-TC dataset, different exams may have different spacings/resolutions, which is the actual distance between two adjacent pixels in a slice. There are various resolutions ranging from 320×320 to 640×640 , and various intervals ranging from $4.4mm$ to $4.6mm$. We resize all slices to 640×640 pixels after spacing alignment to $0.4375mm \times 0.4375mm$ before feeding them into the network. Typical data augmentations including random horizontal flip, random zooming in/out ([0.7, 1.3]) and random crop (to the size of 512×512 pixels) are also applied during training. All models are trained using NVIDIA V100 GPUs without additional tricks.

D. Results and Analysis

Overall, SpineOne surpasses both one-stage and two-stage state-of-the-art methods by integrating all three proposed components into existing one-stage, anchor-free detection methods. Different components are able to boost each other to achieve the best performance across extensive one-stage models. SpineOne with DeepLabv3+(Res101) performs the best among all because ASPP is an excellent ingredient to enlarge the reception field, thus helping fuse multi-level features learned from the backbone. Moreover, SpineOne with CenterNet(Res101) also performs competitively well because CenterNet models an object as the centroid of its bounding box, which is a perfect design for our task when all anatomical structures are annotated as centroids. Visualization results are available in Fig. 9. In addition, our framework can be easily generalized to other analogue tasks, especially medical diagnosis applications. For example, the simple yet effective OCPC heatmap can be applied to other localization and classification tasks

TABLE II
PERFORMANCE COMPARISON OF OUR FRAMEWORK WITH THE TOP THREE WINNING TEAMS IN SDID-TC. THE SCORE IS THE MICRO AVERAGE PRECISION OF ALL CATEGORIES WITH PCK UNDER $6mm$.

Team (ranking) Score	SDID-TC			
	3 rd	2 nd	1 st	SpineOne
	0.679	0.690	0.702	0.716

where keypoints are relatively separated in space. Mapping keypoints into one channel can help capture the geometrical and classification correlations among different keypoints, which is particularly useful for tasks that require the exact location and classification of keypoints. Also, the two attention modules can be extended to tasks where different regions of inputs should be treated more adaptively. Furthermore, the gradient-based OA mechanism is also a potential technique to generally boost the learning of main objectives in multi-head models at the later training stage.

Two-stage methods need to train region proposal networks and classifiers in sequence. They are also more time-consuming during inference (e.g., our SpineOne is only 0.029 seconds during inference on average, in contrast, two-stage methods cost more than 0.3 seconds on average). However, it is surprising that two-stage methods do not outperform one-stage ones as empirically demonstrated in Table I. One important reason is that different from two-stage methods, one-stage methods are able to localize and classify all keypoints simultaneously, better leveraging their intrinsic correlations. For example, the degeneration of disc D_2 may also indicate the degeneration of adjacent discs D_1 and D_3 , while both D_4 and D_5 are normal (Fig. 1). As a result, one-stage methods can take advantage of their geometrical and classification correlations as to predict their locations and classes more accurately than two-stage ones, which classify keypoints by cropping their surrounding regions separately. Our SpineOne can further improve the performance of one-stage methods integrating the OCPC keypoint heatmap, attention modules and the gradient-based OA mechanism.

Moreover, we also compare our results with the top three winning teams in SDID-TC [59]. Two of them (the 1st and 3rd teams) converted the keypoint detection task into an object detection task by changing the ground truth annotations. Then, they exploited state-of-the-art detection methods to obtain good performance. The 2nd team developed a one-stage detection framework for simultaneous localization and classification of degenerative discs and vertebrae, which is similar to our method. Table II reports the final score of SpineOne, which is higher than the best result in the competition. The score is the *micro average precision*, defined as $score = \frac{TP}{TP+FP}$ of all categories with PCK under $6mm$. Again, we compute the *macro mean* of *precision*, *recall*, and *F1 score* towards overall evaluation of the final performance. These performance metrics allow the evaluation of our SpineOne and state-of-the-art methods in a more comprehensive manner.

TABLE III

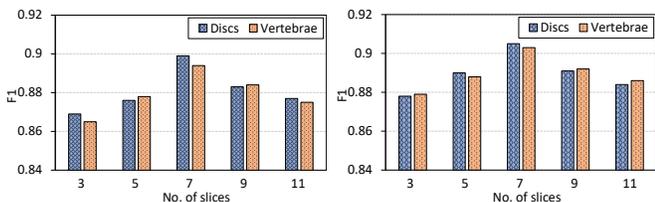
ABLATION STUDY OF OUR FRAMEWORK ON THE SDID-TC DATASET. THE BACKBONE ARCHITECTURES ARE DEEPLABV3+(RES50/101). ALL EVALUATION RESULTS ARE THE MACRO MEAN OF NORMAL AND DEGENERATIVE CLASSES. THE BEST RESULTS IN EVERY COLUMN ARE **BOLDFACED**.

Backbone	Component				Discs			Vertebrae		
	OCPC	PAM	CAM	OA	Recall	Precision	F1	Recall	Precision	F1
DeepLabv3+ (Res50)					0.813	0.848	0.830	0.763	0.784	0.774
	✓				0.827	0.856	0.841	0.783	0.797	0.790
	✓	✓			0.838	0.865	0.851	0.790	0.812	0.801
	✓		✓		0.829	0.859	0.843	0.785	0.802	0.793
	✓	✓	✓		0.841	0.868	0.854	0.800	0.823	0.811
	✓		✓	✓	0.823	0.851	0.836	0.776	0.792	0.784
	✓	✓	✓	✓	0.832	0.861	0.846	0.785	0.806	0.795
	✓	✓	✓	✓	0.847	0.873	0.860	0.811	0.834	0.823
DeepLabv3+ (Res101)					0.819	0.855	0.836	0.774	0.795	0.785
	✓				0.836	0.870	0.852	0.796	0.816	0.806
	✓	✓			0.850	0.882	0.865	0.816	0.838	0.827
	✓		✓		0.847	0.878	0.861	0.810	0.832	0.821
	✓	✓	✓		0.853	0.885	0.868	0.826	0.848	0.837
	✓		✓	✓	0.830	0.864	0.846	0.790	0.810	0.800
	✓	✓	✓	✓	0.844	0.875	0.859	0.805	0.827	0.816
	✓	✓	✓	✓	0.859	0.891	0.874	0.840	0.859	0.849

TABLE IV

PERFORMANCE COMPARISON BETWEEN W/ (EVEN ROWS) AND W/O (ODD ROWS) THE OFFSET MAP AS OUTPUTS.

Backbone	Offset	Component			Discs			Vertebrae		
		OCPC	PAM	CAM	Recall	Precision	F1	Recall	Precision	F1
DeepLabv3+ (Res50)		✓			0.824	0.853	0.838	0.774	0.796	0.785
	✓	✓			0.827	0.856	0.841	0.783	0.797	0.790
		✓	✓	✓	0.837	0.862	0.849	0.792	0.817	0.804
	✓	✓	✓	✓	0.841	0.868	0.854	0.800	0.823	0.811
DeepLabv3+ (Res101)		✓			0.833	0.861	0.846	0.791	0.816	0.803
	✓	✓			0.836	0.870	0.852	0.796	0.816	0.806
		✓	✓	✓	0.849	0.874	0.861	0.820	0.843	0.831
	✓	✓	✓	✓	0.853	0.885	0.868	0.826	0.848	0.837



(a) DeepLabv3+(Res50) (b) DeepLabv3+(Res101)

Fig. 7. Comparison of using different numbers of MRI slices as inputs.

E. Ablation Studies

Table III reports the experimental results of our framework under various settings: 1) different backbones; and 2) with different components (i.e., the OCPC keypoint heatmap, PAM and CAM attention modules, and the gradient-based OA mechanism).

As can be observed, each of our proposed component can improve stage-of-the-art one-stage methods, and our integrated SpineOne framework outperforms them by 3% – 6.4% (abso-

lute) in terms of macro mean F1, to which precision contributes more than recall. Note that the conventional heatmap will have 7 channels (5 centroids and 2 classes) for discs and vertebrae without OCPC, respectively. Our simple OCPC design with only 2 channels (2 classes) can help improve the performance by approximately 2%. PAM performs better than CAM dominantly, which indicates that the inner-image spatial interaction is more important than the inter-channel correlation in our task. The gradient-based OA mechanism can also help improve the final performance of our framework with or without attention modules. Note that the framework with attention modules does not perform well without the OCPC heatmap, which proves the necessity of the OCPC design when diagnosing degenerative discs and vertebrae.

F. More Explorations

1) *Slices selection*: There are several sequences of MRIs in each exam, e.g., sagittal T1 weighted, sagittal T2 weighted, and axial T2 weighted. The cerebrospinal fluid (CSF) is highlighted (brightly shown) in T2 sequences whereas dark

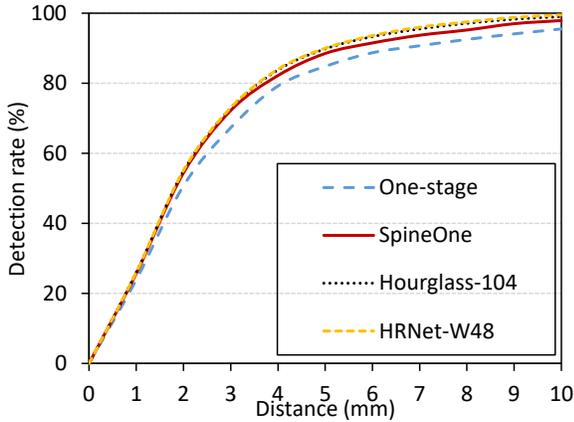


Fig. 8. PCK comparison results. We select DeepLabv3+(Res101) as the backbone for both the conventional one-stage method and SpineOne. We select Hourglass-104 and HRNet-W48 as examples of two-stage methods.

in T1 sequences. Clinical experts mainly focus on the slices in the middle of sagittal T2 to diagnose degenerative discs and vertebrae [6]. As we explained earlier, this is also the reason why we take the slice in the middle of sagittal T2 and its nearby ones as inputs. However, there are usually 11 or more slices in sagittal T2. Therefore, we conduct a comparison study to test the performance on different number of input slices (i.e., slices in number of [3, 5, 7, 9, 11]). This is done with OCPC, but without attention modules or OA for the convenience of comparison and avoiding the effect of other factors. As shown in Fig. 7, the performance is consistent across the two backbones and the best results are obtained when 7 slices in the middle of sagittal T2 are considered. This indicates that slices in the middle of sagittal T2 can indeed help models identify degenerative discs and vertebrae, while those far away from the middle are less useful.

2) *With or without the offset*: Since the offset learning is less important than the keypoint localization and classification, here we further test the performance of our framework with or without the two (i.e., discs and vertebrae) offset outputs, with results reported in Table IV. The short-range offset map uses pixels around the keypoint as the supplementary information to locate it. Note that our OA mechanism has to be removed in this case as it is designed to better associate the offset with the heatmap. It can be observed that models with the offset map (even rows) can generally perform better than those without the offset map (odd rows), to limited extent (less than 1% improvement on the F1 score in Table IV). As empirically proved by the results in Table III, when the offset is considered, our OA mechanism can further improve the performance by approximately 1% on average.

3) *Detection rate*: Although the localization accuracy is not our final goal, it is still worth inspecting it to better understand our framework. Here, we conduct the last experiment to evaluate models by only the detection/identification rate of all centroids, which is calculated by the standard Percentage of Correct Keypoints (PCK) metric in keypoint detection tasks.

TABLE V
PERFORMANCE COMPARISON OF OUR FRAMEWORK WITH EXISTING BASELINES IN THE MICCAI 2014 COMPUTATIONAL CHALLENGE ON VERTEBRAE LOCALIZATION AND IDENTIFICATION.

Independent test	Chen et al. [14]	Yang et al. [60]	Liao et al. [61]	Chen et al. [62]	Ours
Detection rate	84.16%	85.00%	88.30%	94.67%	94.55%

PCK can be formulated in our task as follows:

$$PCK = \frac{TP + FP}{TP + FP + FN}. \quad (8)$$

PCK under $6mm$ (PCK-6) was introduced in Section IV-B suggested by experts. Here, we calculate the detection rate by PCK in an extensive range from $1mm$ to $10mm$ (Fig. 8). We compare four models, i.e., the one-stage method with DeepLabv3+(Res101), SpineOne with DeepLabv3+(Res101), Hourglass-104, and HRNet-W48, which are 2^{nd} , 4^{th} , 9^{th} , and 15^{th} models in Table I. Although there are some variations, our proposed framework demonstrates an approximately 2 – 3% improvement over the one-stage method across the extensive range. For example, PCK-6 of the one-stage method is improved from 88.7% to 91.5%, while PCK-10 is improved from 95.5% to 97.9% by our SpineOne. We also observe that the improvement of the detection rate evaluated by PCK is less significant than those results reported in Table I. The reason is that the proposed OA mechanism encourages our framework to focus on the more important keypoint learning task, rather than the offset learning task. Indeed, the pathological classification in the keypoint learning task plays a more important role in spinal diagnosis than only finding the optimal location of the discs or vertebrae by offset learning. We could actually tolerate some small errors in centroids localization, which is why PCK-6 was set by experts. Interestingly, two-stage methods outperform SpineOne in keypoint localization by about 1% across the extensive range. However, it drops much in classification as it cropped the entire input into small region proposals based on centroids’ locations, which fails to capture the classification correlation among discs and vertebrae.

In addition, we also compare our method with existing baselines on a public dataset from the MICCAI 2014 Computational Challenge on Vertebrae Localization and Identification (Table V), which is not exactly the same task but close to ours. The dataset contains 242 spine-focused CT training scans of various types of high-grade pathologies and metal implants, together with 60 other scans for hold-out evaluation. The differences between this task and our task are two-fold: 1) the input images are 3D CT volumes in this task while the inputs are 2D MRI slices in our task; 2) there is no degeneration classification in this task, so we only use the detection rate as the evaluation metric. We modify the structure of our framework so that it can take as input the 3D CT volumes ($512 \times 512 \times 128$). We empirically show that our method achieves the state-of-the-art performance on this public dataset, where the baseline results are directly copied from their original papers.

TABLE VI
STRUCTURE OF OUR FRAMEWORK. CONV2D() DENOTES CONV2D(IN_CHANNEL, OUT_CHANNEL, KERNEL, STRIDE). BN DENOTES BATCH NORMALIZATION.

Modules	Parameters			
Backbone	e.g., DeepLabv3+			
Two substreams	Discs		Vertebrae	
Attention	PAM+CAM		PAM+CAM	
CNN	Conv2d(64, 32, 3, 1)	Conv2d(64, 32, 3, 1)	Conv2d(64, 32, 3, 1)	Conv2d(64, 32, 3, 1)
	BN	BN	BN	BN
	Conv2d(32, 32, 3, 1)	Conv2d(32, 32, 3, 1)	Conv2d(32, 32, 3, 1)	Conv2d(32, 32, 3, 1)
	BN	BN	BN	BN
Generation	Conv2d(32, 2, 1, 1)	Conv2d(32, 4, 1, 1)	Conv2d(32, 2, 1, 1)	Conv2d(32, 4, 1, 1)
Output	Heatmap \hat{h}_{disc}	Offset \hat{o}_{disc}	Heatmap \hat{h}_{vert}	Offset \hat{o}_{vert}

Ground truth of discs		Generation of discs	
Disc	Coordinate	Class	Coordinate
D ₁	[288, 235]	N	[290, 234]
D ₂	[278, 293]	N	[280, 291]
D ₃	[272, 352]	N	[274, 350]
D ₄	[274, 410]	D	[276, 410]
D ₅	[284, 465]	D	[288, 464]

Ground truth of vertebrae		Generation of vertebrae	
Vertebra	Coordinate	Class	Coordinate
V ₁	[292, 208]	N	[295, 210]
V ₂	[284, 263]	D	[284, 264]
V ₃	[274, 322]	N	[274, 323]
V ₄	[272, 382]	N	[272, 383]
V ₅	[278, 439]	D	[281, 439]

Ground truth of discs		Generation of discs	
Disc	Coordinate	Class	Coordinate
D ₁	[263, 144]	N	[261, 143]
D ₂	[257, 194]	N	[256, 191]
D ₃	[256, 242]	D	[252, 240]
D ₄	[261, 294]	N	[258, 291]
D ₅	[276, 343]	N	[271, 341]

Ground truth of vertebrae		Generation of vertebrae	
Vertebra	Coordinate	Class	Coordinate
V ₁	[263, 122]	D	[263, 120]
V ₂	[257, 170]	D	[258, 167]
V ₃	[252, 218]	D	[252, 216]
V ₄	[253, 266]	D	[253, 267]
V ₅	[262, 318]	D	[260, 319]

Ground truth of discs		Generation of discs	
Disc	Coordinate	Class	Coordinate
D ₁	[301, 197]	N	[302, 195]
D ₂	[283, 246]	D	[280, 244]
D ₃	[266, 300]	D	[267, 298]
D ₄	[266, 359]	D	[266, 355]
D ₅	[289, 407]	D	[289, 403]

Ground truth of vertebrae		Generation of vertebrae	
Vertebra	Coordinate	Class	Coordinate
V ₁	[309, 173]	N	[309, 174]
V ₂	[289, 222]	D	[290, 220]
V ₃	[271, 274]	D	[274, 272]
V ₄	[266, 330]	D	[265, 328]
V ₅	[275, 385]	D	[275, 384]

Ground truth of discs		Generation of discs	
Disc	Coordinate	Class	Coordinate
D ₁	[183, 112]	N	[178, 111]
D ₂	[172, 143]	N	[174, 144]
D ₃	[168, 178]	N	[168, 178]
D ₄	[169, 214]	D	[162, 215]
D ₅	[182, 248]	D	[175, 239]

Ground truth of vertebrae		Generation of vertebrae	
Vertebra	Coordinate	Class	Coordinate
V ₁	[188, 97]	N	[178, 95]
V ₂	[178, 127]	N	[174, 128]
V ₃	[169, 160]	N	[170, 161]
V ₄	[168, 196]	D	[165, 197]
V ₅	[172, 232]	D	[165, 227]

Fig. 9. Random inference examples of generated results by our framework, compared with the ground truth. In each case, the left image is the middle slice with generated (red) and ground truth centroids of discs (green) and vertebrae (blue). The middle/right table shows ground truth/generated pixel locations and degenerative classes. The first three examples are successful cases while the last (bottom right) one is a failure case. Our approach demonstrates good performance in most cases when input slices of the sagittal T2 sequence are with major spacings among the dataset. However, when inputs are with few spacings among the dataset (e.g., slices in the last case are with the extremely large spacing, namely the low resolution), our framework may detect centroids not accurately or misclassify them into wrong classes. The black results in the four tables on the right are good results while the marked results in red are bad ones.

V. CONCLUSION

We propose a one-stage detector to simultaneously localize and classify degenerative discs and vertebrae on the lumbar spine. It addresses two essential tasks of the diagnosis process: 1) anatomical localization; and 2) pathological classification. SpineOne is built upon CNNs with three novel and generic techniques. First, a new design of the OCPC heatmap facilitates above two tasks at the same time. OCPC can help capture both geometrical and classification correlations among keypoints, and can be easily applied to other medical domains where keypoints are spatially separated. Second, dual self-attention modules better differentiate the learning stream for discs and vertebrae, which can be readily plugged into other frameworks to improve discriminative representation learn-

ing. Third, a novel gradient-guided OA mechanism further associates different learning objectives. It explicitly connects heatmap learning with offset learning, thus making the learning of the main objective more effective. We recommend using OA after 75% training epochs when models start to converge. Experimental results show that we can improve state-of-the-art methods using each individual technique and the integrated framework wins by a large margin on the SDID-TC dataset. For future work, we plan to 1) augment the volume of the SDID-TC dataset; 2) research on the degenerative grading for discs and vertebrae; and 3) verify our three proposed techniques in other medical scenarios.

ACKNOWLEDGMENT

We would like to thank clinical experts who have helped build the Spinal Disease Intelligent Diagnosis Tianchi Competition (SDID-TC) dataset. We are grateful for anonymous patients who have provided their MRI exams for our research. We are also thankful for the technical support from our colleagues at the DAMO Academy, Alibaba Group.

REFERENCES

- [1] UCD, *UC Davis Spine Center*, accessed 2020. [Online]. Available: <https://health.ucdavis.edu/spine/specialties/degenerative.html>
- [2] NY, *The Spine Hospital at the Neurological Institute of New York*, accessed 2020. [Online]. Available: <https://www.columbiaspine.org/condition/degenerative-spine-conditions/>
- [3] S. Kadoury, H. Labelle, and N. Paragios, "Spine segmentation in medical images using manifold embeddings and higher-order mrf's," *IEEE Trans. on Medical Imaging*, vol. 32, no. 7, pp. 1227–1238, 2013.
- [4] C. Chen, D. Belavy, W. Yu, C. Chu, G. Armbrrecht, M. Bansmann, D. Felsenberg, and G. Zheng, "Localization and segmentation of 3d intervertebral discs in mr images by data driven estimation," *IEEE Trans. on Medical Imaging*, vol. 34, no. 8, pp. 1719–1729, 2015.
- [5] R. Korez, B. Ibragimov, B. Likar, F. Pernuš, and T. Vrtovec, "A framework for automated spine and vertebrae interpolation-based detection and model-based segmentation," *IEEE Trans. on Medical Imaging*, vol. 34, no. 8, pp. 1649–1662, 2015.
- [6] A. Jamaludin, T. Kadir, and A. Zisserman, "Spinenet: automatically pinpointing classification evidence in spinal mris," in *MICCAI*. Springer, 2016, pp. 166–175.
- [7] Z. Han, B. Wei, A. Mercado, S. Leung, and S. Li, "Spine-gan: Semantic segmentation of multiple spinal structures," *Medical Image Analysis*, vol. 50, pp. 23–35, 2018.
- [8] M. Levine, T. De Silva, M. D. Ketcha, R. Vijayan, S. Doerr, A. Uneri, S. Vedula, N. Theodore, and J. H. Siewerdsen, "Automatic vertebrae localization in spine ct: a deep-learning approach for image guidance and surgical data science," in *Medical Imaging 2019*, vol. 10951. ISOP, 2019, p. 109510S.
- [9] Aliyun, *Tianchi*, accessed 2020. [Online]. Available: <https://tianchi.aliyun.com/competition/entrance/531796/information>
- [10] A. Sekuboyina, J. Kukačka, J. S. Kirschke, B. H. Menze, and A. Valentinič, "Attention-driven deep learning for pathological spine segmentation," in *International Workshop and Challenge on Computational Methods and Clinical Applications in Musculoskeletal Imaging*. Springer, 2017, pp. 108–119.
- [11] G. Zheng, C. Chu, D. L. Belavý, B. Ibragimov, R. Korez, T. Vrtovec, H. Hutt, R. Everson, J. Meakin, I. L. Andrade *et al.*, "Evaluation and comparison of 3d intervertebral disc localization and segmentation methods for 3d t2 mr data: A grand challenge," *Medical image analysis*, vol. 35, pp. 327–344, 2017.
- [12] S. Pang, C. Pang, L. Zhao, Y. Chen, Z. Su, Y. Zhou, M. Huang, W. Yang, H. Lu, and Q. Feng, "Spineparsenet: Spine parsing for volumetric mr image by a two-stage segmentation framework with semantic image representation," *IEEE Trans. on Medical Imaging*, 2020.
- [13] B. Glocker, J. Feulner, A. Criminisi, D. R. Haynor, and E. Konukoglu, "Automatic localization and identification of vertebrae in arbitrary field-of-view ct scans," in *MICCAI*. Springer, 2012, pp. 590–598.
- [14] H. Chen, C. Shen, J. Qin, D. Ni, L. Shi, J. C. Cheng, and P.-A. Heng, "Automatic localization and identification of vertebrae in spine ct via a joint learning model with deep neural networks," in *MICCAI*. Springer, 2015, pp. 515–522.
- [15] S. Zhao, X. Wu, B. Chen, and S. Li, "Automatic vertebrae recognition from arbitrary spine mri images by a category-consistent self-calibration detection framework," *Medical Image Analysis*, vol. 67, p. 101826, 2021.
- [16] J. Chmelik, R. Jakubicek, P. Walek, J. Jan, P. Ourednicek, L. Lambert, E. Amadori, and G. Gavelli, "Deep convolutional neural network-based segmentation and classification of difficult to define metastatic spinal lesions in 3d ct data," *Medical Image Analysis*, vol. 49, pp. 76–88, 2018.
- [17] J.-T. Lu, S. Pedemonte, B. Bizzo, S. Doyle, K. P. Andriole, M. H. Michalski, R. G. Gonzalez, and S. R. Pomerantz, "Deepspine: Automated lumbar vertebral segmentation, disc-level designation, and spinal stenosis grading using deep learning," *arXiv preprint*, vol. arXiv, p. 21807.10215, 2018.
- [18] D. Zhang, B. Chen, and S. Li, "Sequential conditional reinforcement learning for simultaneous vertebral body detection and segmentation with modeling the spine anatomy," *Medical Image Analysis*, vol. 67, p. 101861, 2021.
- [19] M. Rashid, X. Gu, and Y. Jae Lee, "Interspecies knowledge transfer for facial keypoint detection," in *CVPR*, 2017, pp. 6894–6903.
- [20] T. Simon, H. Joo, I. Matthews, and Y. Sheikh, "Hand keypoint detection in single images using multiview bootstrapping," in *CVPR*, 2017, pp. 1145–1153.
- [21] G. Georgakis, S. Karanam, Z. Wu, J. Ernst, and J. Košecká, "End-to-end learning of keypoint detector and descriptor for pose invariant 3d matching," in *CVPR*, 2018, pp. 1965–1973.
- [22] A. Toshev and C. Szegedy, "DeepPose: Human pose estimation via deep neural networks," in *CVPR*, 2014, pp. 1653–1660.
- [23] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler, "Joint training of a convolutional network and a graphical model for human pose estimation," in *NeurIPS*, 2014, pp. 1799–1807.
- [24] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," in *CVPR*, 2016, pp. 4724–4732.
- [25] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *CVPR*, 2017, pp. 7291–7299.
- [26] G. Hidalgo, Y. Raaj, H. Idrees, D. Xiang, H. Joo, T. Simon, and Y. Sheikh, "Single-network whole-body pose estimation," in *ICCV*, 2019, pp. 6982–6991.
- [27] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu, "Rmpe: Regional multi-person pose estimation," in *ICCV*, 2017, pp. 2334–2343.
- [28] Y. Chen, Z. Wang, Y. Peng, Z. Zhang, G. Yu, and J. Sun, "Cascaded pyramid network for multi-person pose estimation," in *CVPR*, 2018, pp. 7103–7112.
- [29] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *ECCV*. Springer, 2016, pp. 483–499.
- [30] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *CVPR*, 2019, pp. 5693–5703.
- [31] W. Yang, S. Li, W. Ouyang, H. Li, and X. Wang, "Learning feature pyramids for human pose estimation," in *ICCV*, 2017, pp. 1281–1290.
- [32] X. Chu, W. Yang, W. Ouyang, C. Ma, A. L. Yuille, and X. Wang, "Multi-context attention for human pose estimation," in *CVPR*, 2017, pp. 1831–1840.
- [33] G. Papandreou, T. Zhu, L.-C. Chen, S. Gidaris, J. Tompson, and K. Murphy, "Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model," in *ECCV*, 2018, pp. 269–286.
- [34] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *ECCV*, 2018, pp. 801–818.
- [35] T.-J. Yang, M. D. Collins, Y. Zhu, J.-J. Hwang, T. Liu, X. Zhang, V. Sze, G. Papandreou, and L.-C. Chen, "Deeperlabs: Single-shot image parser," *arXiv preprint*, vol. arXiv, p. 1902.05093, 2019.
- [36] H. Law and J. Deng, "Cornernet: Detecting objects as paired keypoints," in *ECCV*, 2018, pp. 734–750.
- [37] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "Centernet: Keypoint triplets for object detection," in *ICCV*, 2019, pp. 6569–6578.
- [38] X. Zhou, J. Zhuo, and P. Krahenbuhl, "Bottom-up object detection by grouping extreme and center points," in *CVPR*, 2019, pp. 850–859.
- [39] Z. Dong, G. Li, Y. Liao, F. Wang, P. Ren, and C. Qian, "Centripetalnet: Pursuing high-quality keypoint pairs for object detection," in *CVPR*, 2020, pp. 10519–10528.
- [40] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv preprint*, vol. arXiv, p. 1904.07850, 2019.
- [41] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *ICCV*, 2019, pp. 9627–9636.
- [42] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [43] J. He, S. Erfani, S. Wijewickrema, S. O'Leary, and K. Ramamohanarao, "Learning non-unique segmentation with reward-penalty dice loss," in *IJCNN*. IEEE, 2020, pp. 1–8.

- [44] J. He, S. Erfani, X. Ma, J. Bailey, Y. Chi, and X.-S. Hua, "Alpha-iou: A family of power intersection over union losses for bounding box regression," *arXiv preprint arXiv:2110.13675*, 2021.
- [45] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *NeurIPS*, 2017, pp. 5998–6008.
- [46] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *CVPR*, 2018, pp. 7794–7803.
- [47] S. Woo, J. Park, J.-Y. Lee, and I. So Kweon, "Cbam: Convolutional block attention module," in *ECCV*, 2018, pp. 3–19.
- [48] S. Liu, E. Johns, and A. J. Davison, "End-to-end multi-task learning with attention," in *CVPR*, 2019, pp. 1871–1880.
- [49] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, "Self-attention generative adversarial networks," in *ICML*, 2019, pp. 7354–7363.
- [50] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in *CVPR*, 2019, pp. 3146–3154.
- [51] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *ICCV*, 2017, pp. 2961–2969.
- [52] A. Newell, Z. Huang, and J. Deng, "Associative embedding: End-to-end learning for joint detection and grouping," in *NeurIPS*, 2017, pp. 2277–2287.
- [53] G. Papandreou, T. Zhu, N. Kanazawa, A. Toshev, J. Tompson, C. Bregler, and K. Murphy, "Towards accurate multi-person pose estimation in the wild," in *CVPR*, 2017, pp. 4903–4911.
- [54] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017, pp. 2980–2988.
- [55] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [56] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *CVPR*, 2017, pp. 4700–4708.
- [57] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017, pp. 2117–2125.
- [58] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint*, vol. arXiv, p. 1412.6980, 2014.
- [59] Aliyun, *RankingList*, accessed 2020. [Online]. Available: <https://tianchi.aliyun.com/competition/entrance/531796/rankingList/2>
- [60] D. Yang, T. Xiong, D. Xu, Q. Huang, D. Liu, S. K. Zhou, Z. Xu, J. Park, M. Chen, T. D. Tran *et al.*, "Automatic vertebra labeling in large-scale 3d ct using deep image-to-image network with message passing and sparsity regularization," in *International conference on information processing in medical imaging*. Springer, 2017, pp. 633–644.
- [61] H. Liao, A. Mesfin, and J. Luo, "Joint vertebrae identification and localization in spinal ct images by combining short-and long-range contextual information," *IEEE Trans. on Medical Imaging*, vol. 37, no. 5, pp. 1266–1275, 2018.
- [62] Y. Chen, Y. Gao, K. Li, L. Zhao, and J. Zhao, "vertebrae identification and localization utilizing fully convolutional networks and a hidden markov model," *IEEE Trans. on Medical Imaging*, vol. 39, no. 2, pp. 387–399, 2019.