

HCA-NET: HIERARCHICAL CONTEXT ATTENTION NETWORK FOR INTERVERTEBRAL DISC SEMANTIC LABELING

Afshin Bozorgpour¹ Bobby Azad² Reza Azad³ Yury Velichko⁴ Ulas Bagci⁴ Dorit Merhof^{1,5}

¹Faculty of Informatics and Data Science, University of Regensburg, Germany

²South Dakota State University, Brookings, USA

³Faculty of Electrical Engineering and Information Technology, RWTH Aachen University, Germany

⁴Machine and Hybrid Intelligence Lab, Northwestern University, Chicago, IL, USA

⁵Fraunhofer Institute for Digital Medicine MEVIS, Germany

ABSTRACT

Accurate and automated segmentation of intervertebral discs (IVDs) in medical images is crucial for assessing spine-related disorders, such as osteoporosis, vertebral fractures, or IVD herniation. We present HCA-Net, a novel contextual attention network architecture for semantic labeling of IVDs, with a special focus on exploiting prior geometric information. Our approach excels at processing features across different scales and effectively consolidating them to capture the intricate spatial relationships within the spinal cord. To achieve this, HCA-Net models IVD labeling as a pose estimation problem, aiming to minimize the discrepancy between each predicted IVD location and its corresponding actual joint location. In addition, we introduce a skeletal loss term to reinforce the model's geometric dependence on the spine. This loss function is designed to constrain the model's predictions to a range that matches the general structure of the human vertebral skeleton. As a result, the network learns to reduce the occurrence of false predictions and adaptively improves the accuracy of IVD location estimation. Through extensive experimental evaluation on multi-center spine datasets, our approach consistently outperforms previous state-of-the-art methods on both MRI T1w and T2w modalities. The codebase is accessible to the public on [GitHub](#).

Index Terms— Intervertebral Disc, Semantic Labeling.

1. INTRODUCTION

The human spinal column comprises 33 individual vertebrae, organized in a stacked configuration and interconnected by ligaments and intervertebral discs (IVDs), commonly referred to as IVDs. This anatomical structure is further categorized into five distinct regions, including the cervical, thoracic, lumbar, sacral, and caudal vertebrae [1]. Each of these regions plays a critical role in various physiological functions, such as shock absorption, load bearing, spinal cord protection, and load distribution management [2]. IVDs are fibrocartilaginous cushions that serve as primary articulations between

adjacent vertebrae. They play a critical role in absorbing the forces and shocks exerted on the body during movement, ensuring spinal flexibility while preventing vertebral friction. Any disruption to the structural integrity of IVDs, whether due to aging, degeneration, or injury, can alter the properties and affect the mechanical performance of the surrounding tissues. Consequently, the precise localization and segmentation of IVDs are critical steps in the diagnosis of spinal disorders and provide invaluable insights into the efficacy of treatment modalities. To address this challenge, numerous semi-automated and fully automated methods have been proposed in the literature [3, 4, 5, 6].

Gros et al. [7] proposed a local descriptor-based method to detect the C2/C3 intervertebral disc (IVD) in medical imaging. This technique compares the mutual information between a patient's image and a template to find the region closest to the spine template. This handcrafted approach generally yields good results, but its performance degrades significantly when the patient's images deviate significantly from the template. To overcome these limitations of manual methods, deep learning models have been employed for robust IVD labeling. Chen et al. [8] introduced a 3D CNN model for MRI data to enabling 3D segmentation and accurate identification of vertebral disc locations. Cai et al. [9] utilized a 3D Deformable Hierarchical Model for 3D spatial vertebral disc localization. Rouhier et al. [4] trained a Count-ception model on 2D MRI sagittal slices to detect vertebral discs. Adibatti et al. [3] proposed a capsule stacked autoencoder for IVD segmentation. Vania et al. [10] introduced a multi-optimization training system at various stages to enhance computational efficiency, building upon Mask R-CNN. Meanwhile, Wimmer et al. [11] presented a cross-modality method for detecting both vertebral and intervertebral discs in volumetric data, using a local entropy-based texture model followed by alignment and refinement techniques. Mbarki et al. [12] employed transfer learning to detect lumbar discs from axial images using a 2D convolutional structure. Their network, based on the U-Net structure with a VGG backbone, generated a spine segmenta-

tion mask used to calculate herniation in lumbar discs. Azad et al. [13] redefined semantic vertebral disc labeling as pose estimation by implementing an hourglass neural network for semantic labeling of IVDs. In a more recent approach [5], they propose an enhancement to the detection process by including the image gradient as an auxiliary input to better capture and represent global shape information.

Existing methods have attempted to improve shape information by incorporating image gradients as auxiliary data [5], focusing on vertebral column region detection [10], and modeling pose information [13]. However, these methods still face limitations in implicitly conditioning the representation space using global vertebral column information to efficiently model geometric constraints. As a result, these strategies may lead to undesirable false positive and false negative predictions. To address this challenge, we present HCA-Net, a novel pose estimation approach that leverages a robust framework featuring Multi-scale Large Kernel Attention (M-LKA) modules to facilitate the comprehensive capture of contextual information while preserving local intricacies. This architectural enhancement plays a pivotal role in enabling precise

semantic labeling. Furthermore, to enhance the model’s reliance on vertebral column geometry, we introduce the skeleton loss function to effectively constrains the model’s predictions within a range consistent with the human vertebral skeleton. Our key contributions are: (1) A contextual attention network for semantic labeling, which incorporates the multi-scale large kernel attention mechanism to model both local and global representations, (2) the skeleton loss function to implicitly enforce geometrical information of the vertebral column into the model prediction.

2. METHOD

The design of our contextual attention network for IVD labeling is driven by the need to extract information from medical images at different scales. While local features are essential for discerning specific anatomical structures such as IVDs, achieving precise disc labeling requires a holistic understanding of the entire spinal structure. This includes considerations such as the orientation of the spine, the arrangement of the IVDs, and the relationships between neighboring discs, which are most effectively captured at different scales within the medical image. To address this challenge, we introduce our novel hierarchical context attention strategy, illustrated in Figure 1. Our approach incorporates multi-scale, large kernel attention blocks to capture both local and global dependencies, while constraining the model prediction with prior information on the distribution of the IVDs.

2.1. Network Architecture

The architecture of the HCA-Net is structured as follows: First, a sequence of convolutional layers is applied to process the input MRI image and transform it into a latent representation. Next, a hierarchical context attention module is employed to capture multi-scale representations. This module uses an hourglass block [14] to effectively model local representations, and then leverages large kernel attention across multiple scales to adjust the representation space based on local-to-global information, facilitating the incorporation of both local and long-range dependencies.

Figure 1 illustrates the construction of HCA-Net, which involves stacking hierarchical context attention (HCA) blocks and incorporates the process of learning object pose estimation through $(N-1)$ intermediate predictions Out_j along with one final prediction. This approach takes into account the multilevel representations generated by the N -stacked HCA blocks. Finally, we merge the intermediate and final prediction masks using the 1×1 convolution, resulting in a V channel prediction map (\hat{y}). Each channel within this map corresponds to a specific intervertebral location, thus providing a comprehensive representation of intervertebral positions. To minimize the network’s prediction error, we take the sum of

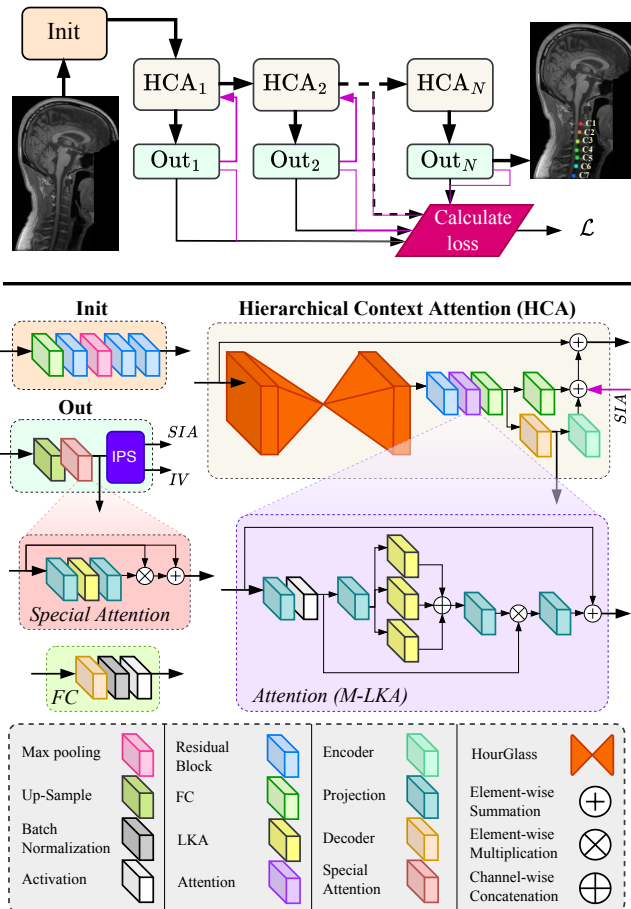


Fig. 1: Structure of the proposed HCA-Net method for IVD semantic labeling.

mean squared error (MSE) loss between the network prediction \hat{y} and the ground truth y :

$$\mathcal{L}_v = \frac{1}{V \times M} \sum_{i=1}^V \sum_{p=1}^M (y_p^i - \hat{y}_p^i)^2, \quad (1)$$

where M corresponds to the number of pixels in the ground truth mask. To reinforce the incorporation of vertebral column structure as an additional supervisory signal to enhance network predictions, we introduce the ‘‘skeleton loss’’ \mathcal{L}_{sk} term to the overall loss function. Consequently, during each training step, HCA-Net aims to minimize the combined loss function:

$$\mathcal{L} = \mathcal{L}_v + \lambda \mathcal{L}_{sk} \quad (2)$$

2.1.1. Multi-scale Large Kernel Attention (M-LKA)

Achieving accurate semantic labeling of IVDs requires the consideration of both local and global semantic representations. Given the geometrical interdependencies among intervertebral joint locations, relying solely on local representations may result in erroneous predictions. To overcome these challenges, we introduce an innovative approach that leverages the Large Kernel Attention (LKA) mechanism. We enhance the LKA module by extending it across multiple scales. The rationale behind this enhancement is to efficiently capture and integrate information at various spatial resolutions, which is especially valuable for tasks demanding precise predictions. In contrast to the original LKA, which employs fixed-sized filters and faces challenges in fully capturing information at different scales within an image, our M-LKA module utilizes parallel filters of varying sizes. This approach allows us to capture both fine-grained details and high-level semantic information concurrently.

The LKA module decomposes a $C \times C$ convolution into three components: a $\lceil \frac{C}{d} \rceil \times \lceil \frac{C}{d} \rceil$ depth-wise dilation convolution (*DW-D-Conv*) for long-range spatial convolution, a $(2d - 1) \times (2d - 1)$ depth-wise convolution (*DW-Conv*) for local spatial convolution, and a 1×1 convolution for channel-wise convolution. This decomposition enables us to extract long-range relationships within the feature space while maintaining computational efficiency and a manageable parameter count when generating the attention map. We further extend the LKA module into multiscale form as follows:

Let $F_S(x)$ represent a set of feature maps obtained by applying depth-wise convolution (*DW-Conv*) to the input features $F(x)$ for each scale $s \in \mathbb{S}$. Then, $F_S(x)$ can be expressed as:

$$F_S(x) = \{(\text{DW-Conv}(F(x)))_s \mid s \in \mathbb{S}\}$$

Subsequently, the attention map *Attention* is generated by applying a 1×1 convolution ($\text{Conv}_{1 \times 1}$) to the feature maps obtained through depth-wise dilation convolution

(*DW-D-Conv*) of $F_S(x)$:

$$\text{Attention} = \text{Conv}_{1 \times 1}(\text{DW-D-Conv}(F_S(x)))$$

Finally, the output x' is computed as the element-wise multiplication (\otimes) between the attention map *Attention* and the input features $F(x)$:

$$x' = \text{Attention} \otimes F(x)$$

2.2. Skeleton Loss Function

Accurate IVD semantic labeling often faces the challenge of generating false predictions, necessitating a mechanism for guiding the network towards more reliable outcomes. To tackle this issue, we leverage the network’s prediction map, denoted as \hat{y} , and apply the softmax operation to transform it into a 2D positional probability distribution for the IVD location in each channel:

$$\mathbf{P}_j^i = \frac{\sigma(\hat{y}^i)}{\sum_p^M \sigma(\hat{y}^i)_p}, \quad (3)$$

where \mathbf{P}_j^i represents the probability of the respective intervertebral joint location within each channel. Subsequently, we use the probability map to generate prototypes for each intervertebral location through T times sampling from each channel and averaging as follows:

$$\mathbf{V}_j^i = \frac{1}{T} \sum_T \text{Sampler}(\mathbf{P}_j^i),$$

The *sampler* function utilizes the probability map \mathbf{P}_j^i to extract intervertebral locations in each channel. Subsequently, our approach integrates a distance function denoted as $D : \mathbb{R}^M \times \mathbb{R}^M \rightarrow [0, +\infty)$ to minimize the distance between the intervertebral column and the ground truth location. To this end, we model the skeleton loss function as follows:

$$\begin{aligned} \mathcal{L}_{sk} &= \sum_{j=1}^N \left(\beta \mathcal{L}_j^{id} + (1 - \beta) \mathcal{L}_j^{pd} \right) \\ \mathcal{L}_j^{id} &= \|\mathbf{V}_j^i - \mathbf{V}_j^{\text{GT}}\|, \quad \mathcal{L}_j^{pd} = \text{PD}(\mathbf{V}_j^i, \mathbf{V}_j^{\text{GT}}) \\ \text{PD}(\mathbf{V}, \mathbf{V}^{\text{GT}}) &= \sum_c^{C-1} \sum_{k=c}^C \alpha^{k-c} (D(\mathbf{V}, c, k) - D(\mathbf{V}^{\text{GT}}, c, k))^2 \end{aligned}$$

Here, we define the distance function as $D(V, i, k) = \|V_i - V_{i+k}\|$. The parameter α represents a learnable weight, while \mathcal{L}^{id} denotes the L_2 distance between the vertebral column prototype and the ground truth. Additionally, \mathcal{L}^{pd} quantifies the pair-wise distance (PD), ensuring the preservation of the geometrical relationships within the intervertebral skeleton structure.

Table 1: Intervertebral disc semantic labeling on the spine generic public dataset. Note that **DTT** indicates Distance to target

Method	T1			T2		
	DTT (mm)	FNR (%)	FPR (%)	DTT (mm)	FNR (%)	FPR (%)
Template Matching [15]	1.97(±4.08)	8.1	2.53	2.05(±3.21)	11.1	2.11
Countception [4]	1.03(±2.81)	4.24	0.9	1.78(±2.64)	3.88	1.5
Pose Estimation [13]	1.32(±1.33)	0.32	0.0	1.31(±2.79)	1.2	0.6
Look Once	1.2(±1.90)	0.7	0.0	1.28(±2.61)	0.9	0.0
HCA-Net without \mathcal{L}_{sk}	1.27(±1.78)	0.6	0.0	1.34(±2.28)	1.2	0.0
HCA-Net	1.19(±1.08)	0.3	0.0	1.26(±2.16)	0.61	0.0

3. EXPERIMENTAL SETUP AND RESULTS

Experimental Setup: In our experiment, we use the Spine Generic Dataset [16] for IVD labeling. This dataset contains samples from 42 medical centers around the world in both T1-weighted (T1w) and T2-weighted (T2w) contrasts and exhibits a large variation in terms of quality, scale, and imaging device. To prepare the dataset for the training, we first calculate the average of six sagittal slices, centered on the middle slice, to create a representative data sample for each subject. To ensure uniformity and to minimize the impact of data variations, we normalize each image to the [0, 1] range. Next, using the IVD coordinate on the 2D position, we create a heatmap image by applying a Gaussian kernel convolution on each position of the IVD. Similar to [4] we extract 11 IVDs for each subject. In instances where an IVD is missing, we designate its position as “unknown” and mitigate its influence on the training process by effectively filtering it out using the visibility flag within the loss function. Following [13], we train the model for 500 epochs with RMSprop optimization using a learning rate of $2.5e - 4$ and a batch size of 4. Our experimental hyperparameter settings entail $\lambda = 2e - 4$ (in Equation 2), $\beta = 0.75$ (in Equation 2.2) and $\alpha = 0.8$ in the PD function. We follow evaluation metrics from prior studies [13, 4], including L2 distance for predicted vs. ground truth IVD positions in 3D space. Additionally, we report False Positive Rate (FPR) and False Negative Rate (FNR).

Results: Table 1 presents a comprehensive analysis of our HCA-Net compared to other SOTA methods for IVD semantic labeling. Our approach consistently outperforms existing methods in both T1w and T2w MRI modalities, showcasing its superior accuracy and reliability. In T1w MRI, our method excels with an impressive average distance to the target (DTT) of 1.19 mm, significantly outperforming other methods. This low DTT, combined with a standard deviation of only 1.08 mm, makes our approach highly reliable for precise IVD localization. Notably, even without the \mathcal{L}_{sk} module, our HCA-Net performs remarkably well, achieving a DTT of 1.27 mm and displaying superior accuracy compared to the alternatives. In the T2w MRI, our HCA-Net again enhances the performance, with an outstanding DTT of 1.26 mm. This result significantly outperforms previous work, underlining the robustness and accuracy of our approach. Additionally, our method achieves a lower false negative rate (FNR) of 0.61%

in T2w, indicating its ability to capture IVDs effectively and minimize missed detections.

In Figure 2, we provide a visual comparison between our HCA-Net and the pose estimation approach [13] in both T1w and T2w modalities. This comparison highlights the precision of our predictions. While the pose estimation approach misses one intervertebral location in T1w modality, our method successfully recognizes all intervertebral locations, with predictions closely matching the actual locations. This visual demonstration underscores the superior performance and accuracy of our HCA-Net.

Comparing our approach to the alternatives, we observe several key advantages. First, HCA-Net eliminates the need for complex preprocessing steps, such as image straightening or spinal cord region detection used in [4], making it more efficient and user-friendly. Second, our approach takes into account spatial relationships between IVDs, contributing to its superior performance, especially in FNR reduction.

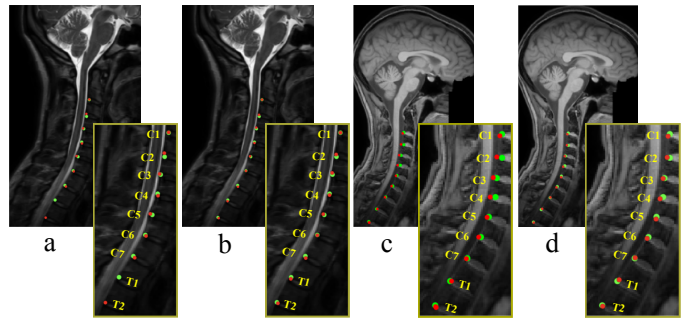


Fig. 2: Comparison of results on T1w (a-b) and T2w (c-d) MRI modalities between the proposed HCA-Net (b and d) and the pose estimation method [13] (a and c). Green dots denote ground truth.

4. CONCLUSION

We proposed HCA-Net, a novel framework that capitalizes on a stack of hierarchical attention blocks to effectively encode both local and global information, ensuring precise localization of IVDs. The incorporation of a skeleton loss function further fine-tunes network predictions by considering the geometry of the intervertebral column. Through comprehensive experimentation, HCA-Net consistently demonstrated superior performance, attaining SOTA results.

5. REFERENCES

- [1] Jan Gewiess, Janick Eglauf, Astrid Soubrier, Sibylle Grad, Mauro Alini, Marianna Peroglio, and Junxuan Ma, “The influence of intervertebral disc overloading on nociceptor calcium flickering,” *JOR Spine*, p. e1267, 2023.
- [2] Ali Al-kubaisi and Nasser N Khamiss, “A transfer learning approach for lumbar spine disc state classification,” *Electronics*, vol. 11, no. 1, pp. 85, 2022.
- [3] Spurthi Adibatti, KR Sudhindra, and Joshi Manisha Shivaram, “Segmentation and classification of intervertebral disc using capsule stacked autoencoder,” *Biomedical Signal Processing and Control*, vol. 86, pp. 105311, 2023.
- [4] Lucas Rouhier, Francisco Perdigon Romero, Joseph Paul Cohen, and Julien Cohen-Adad, “Spine intervertebral disc labeling using a fully convolutional redundant counting model,” *arXiv preprint arXiv:2003.04387*, 2020.
- [5] Reza Azad, Moein Heidari, Julien Cohen-Adad, Ehsan Adeli, and Dorit Merhof, “Intervertebral disc labeling with learning shape information, a look once approach,” *MICCAI 2022*, 2022.
- [6] Chao Hou, Xiaogang Li, Hongbo Wang, Weiqi Zhang, Fei Liu, Defeng Liu, and Yuzhen Pan, “An mri image automatic diagnosis model for lumbar disc herniation using semi-supervised learning,” *Complex & Intelligent Systems*, pp. 1–18, 2023.
- [7] Charley Gros, Benjamin De Leener, Sara M Dupont, Alan R Martin, Michael G Fehlings, Rohit Bakshi, Subhash Tummala, Vincent Auclair, Donald G McLaren, Virginie Callot, et al., “Automatic spinal cord localization, robust to mri contrasts using global curve optimization,” *Medical image analysis*, vol. 44, pp. 215–227, 2018.
- [8] Yizhi Chen, Yunhe Gao, Kang Li, Liang Zhao, and Jun Zhao, “vertebrae identification and localization utilizing fully convolutional networks and a hidden markov model,” *IEEE transactions on medical imaging*, vol. 39, no. 2, pp. 387–399, 2019.
- [9] Yunliang Cai, Said Osman, Manas Sharma, Mark Landis, and Shuo Li, “Multi-modality vertebra recognition in arbitrary views using 3d deformable hierarchical model,” *IEEE transactions on medical imaging*, vol. 34, no. 8, pp. 1676–1693, 2015.
- [10] Malinda Vania and Deukhee Lee, “Intervertebral disc instance segmentation using a multistage optimization mask-rcnn (mom-rcnn),” *Journal of Computational Design and Engineering*, vol. 8, no. 4, pp. 1023–1036, 2021.
- [11] Maria Wimmer, David Major, Alexey A Novikov, and Katja Bühler, “Fully automatic cross-modality localization and labeling of vertebral bodies and intervertebral discs in 3d spinal images,” *International journal of computer assisted radiology and surgery*, vol. 13, no. 10, pp. 1591–1603, 2018.
- [12] Wafa Mbarki, Moez Bouchouicha, Sebastien Frizzi, Frederick Tshibusu, Leila Ben Farhat, and Mounir Sayadi, “Lumbar spine discs classification based on deep convolutional neural networks using axial view mri,” *Interdisciplinary Neurosurgery*, vol. 22, pp. 100837, 2020.
- [13] Reza Azad, Lucas Rouhier, and Julien Cohen-Adad, “Stacked hourglass network with a multi-level attention mechanism: Where to look for intervertebral disc labeling,” in *International Workshop on Machine Learning in Medical Imaging*. Springer, 2021, pp. 406–415.
- [14] Alejandro Newell, Kaiyu Yang, and Jia Deng, “Stacked hourglass networks for human pose estimation,” in *European conference on computer vision*. Springer, 2016, pp. 483–499.
- [15] Eugénie Ullmann, Jean François Pelletier Paquette, William E Thong, and Julien Cohen-Adad, “Automatic labeling of vertebral levels using a robust template-based approach,” *International journal of biomedical imaging*, vol. 2014, 2014.
- [16] J Cohen-Adad and et al, “Open-access quantitative mri data of the spinal cord and reproducibility across participants, sites and manufacturers. sci. data. doi: 10.1038/s41596-021-00588-0,” .