Organs at Risk Segmentation, with a Focus on Spinal Cord, Using Deep Learning Approaches

1

Abstract. Cancer treatment with radiotherapy (RT) presents the challenge of protecting organs adjacent to the treatment area, called Organs at Risk (ORs). This study proposes to automate the segmentation of the spinal cord (SM) and ORs in Computed Tomography (CT) images, essential to aid in radiotherapy planning. The proposed method involves three steps: i) data acquisition, ii) pre-processing and detecting the spine region using YOLOv8, and iii) final segmentation of the ME using the U-Net architecture. The results highlight a Dice precision: 88%, Sensitivity: 92%, Specificity: 99%, Accuracy: 99%. This approach will contribute to more effective radiotherapy planning and improve patient clinical outcomes.

1. Introduction

Organs at risk (ORs) primarily refer to healthy organs near the target volume that may be affected by radiation exposure during radiotherapy (RT) treatment for cancer. These organs [Morbidity 2016] are susceptible to damage during radiotherapy treatment. The treatment is directed at the tumor volume, known as target volume [Shirato and et al 2018].

The radiation dose directed to the target volume aims to eliminate tumor cells, preventing them from increasing. In radiotherapy, high doses of radiation (greater than 50 Gy) are applied to treat malignant tissues, while for benign tissues, lower to intermediate doses (3–50 Gy) are used to control growth [McKeown et al. 2015] effectively. Gray (Gy) is the SI unit of the absorbed radiation dose. A dose of 1 Gray means that one kilogram of tissue has absorbed one joule of radiation energy.

OR segmentation is a crucial step in radiotherapy planning, commonly performed using 3D imaging techniques such as CT. Recently, there has been a significant focus on improving methods for segmenting these at-risk organs, with ongoing research exploring computer vision and deep learning approaches to this end [Lima and Oliveira 2018, Matos et al. 2023, Santos et al. 2023].

Among the ORs, ME stands out, which will be the focus of study in this work. Identifying ORs is fundamental. Therefore, computational methods were developed to assist specialists in segmentation [Diniz et al. 2021a]. Identifying ORs is time-consuming for medical professionals, requiring an extensive team of specialists. The exhaustive analysis of exam by exam becomes tiring and susceptible to errors, especially in organs such as ME, which require segmentation in all slices of CT exams [West et al. 2018].

Therefore, considering the problem presented, this work proposes a method that uses deep learning to perform automated detection and segmentation of ORs, offering a second opinion to experts. As contributions, the following stand out:

- 1. A completely automatic method that surpasses previous results published in the literature.
- 2. A tool to help specialists identify OARs quickly and accurately.

2. Related Work

In this section, works that deal with ME segmentation problems will be presented. Table 1 presents a summary of related work.

Work Dataset and Sample Techniques Used Objective Results Template matching, superpixel, CNN; Adaptive Model Matching, CNN [Diniz et al. 2021b] Spinal cord segmentation on CT AAPM Thoracic Autosegmentation Dice: 78.20% / 81.69% (ME) / 36 CT images [Roncaglioni 2021] OR segmentation in radiotherapy Structseg2019, SegTHOR2019 / 90 Bagging, Boosting, Stacking; Con-Dice: 89.92% (ME) volutional Neural Networks [Lambert et al. 2020] Segmentation of OR's in the chest SegTHOR / 60 CT images U-NET Architectures Segmentation of thoracic OR's on StructSeg 2019 / 50 CT images CT Cascaded SE-ResUnet [Cao et al. 2021] Dice: 91% (ME) Automatic segmentation of OR's in Private dataset / 755 CT images. WBNet [Chen et al. 2021] Dice: 88% (ME) [SPINONI 2021] Automatic segmentation of OR's in StructSeg / 50 CT images. Dice: 84.05% (ME) Convolutional Neural Networks AAPM Thoracic Autosegmentation 3D II-Net [Feng et al. 2019] Segmentation of thoracic ORs on Dice: 89% (ME) / 36 CT images. Private dataset / 16,024 exams. [Vieira et al. 2022] DL-based CAD system Classification Ensemble AUC: 86%, Accuracy: 86%

Table 1. Related Work

2.1. Evaluation of related work

Observing the literature, the following limitations can be highlighted:

- Some studies do not address the segmentation of all relevant organs in certain body regions, such as the cervical spine, which limits the applicability of the methods in broader clinical contexts [Chen et al. 2021].
- Although many methods have shown promising results on specific datasets, there are still challenges in ensuring the generalizability of these methods to different types of exams and clinical conditions [SPINONI 2021].
- Some methods, especially those based on Deep Learning, may require considerable computational resources during training and inference, especially those with 3D approaches [Cao et al. 2021, Feng et al. 2019], which may limit their applicability in clinical environments with limited resources [Roncaglioni 2021].
- Need for Manual Annotations: Most studies still rely on manual annotations from experts to train and evaluate models, which can be time-consuming and susceptible to human error [Lambert et al. 2020].

Thus, the proposed method seeks to address the limitations presented where it is proposed. Deep learning techniques, specifically the You Only Look Once (YOLO) architecture, are recognized for their efficiency in detecting objects in images. By employing YOLOv8 to detect the region of interest, we seek a more comprehensive and accurate approach to identifying organs of interest, such as the ME. Then, applying U-Net to perform fine segmentation of these organs provides a second opinion to experts and contributes to more accurate and effective radiotherapy planning. This integration between automated detection and segmentation has the potential to fill the gaps left by related studies, providing a more complete and robust solution for organ segmentation in medical images.

3. Materials and Method

This section describes the proposed method for segmenting the ME as a risk organ in CT. The method consists of three steps: data acquisition, spine region detection using YOLOv8, and final ME segmentation with the U-Net architecture. Figure 1 illustrates the steps of the proposed method.

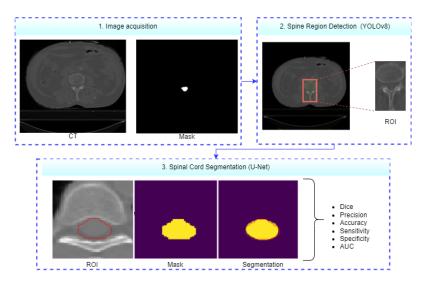


Figure 1. Steps of the proposed method.

3.1. Image Acquisition

The database used in this study originated from the AAPM 2017 Thoracic Autosegmentation Challenge and consists of CT images covering the thoracic region. It comprises 60 exams from different patients from three different institutions, each contributing 20 patients, which gives heterogeneity to the data set. The data was distributed into three groups stratified by institution: 36 training sets, 12 external test sets, and 12 internal test sets. The total set comprises 9,593 slices [Langner et al. 2013]. A 2D approach will be used for this work, using exam slices from these data sets. An example of a database image and mask is illustrated in Figure 2.

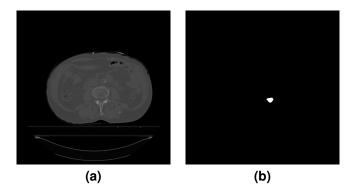


Figure 2. Dataset. (a) Image, (b) Mask.

It is important to highlight that the ME region is extremely small compared to the entire volume. Due to its complexity and anatomical variations, this makes semantic seg-

mentation models less accurate. Therefore, an initial stage of detecting the spine region is necessary to improve the accuracy of the final segmentation.

3.2. Spine region detection

As highlighted, treating the entire volume in the ME segmentation makes the final segmentation work more complex. To facilitate this process, the spine region is detected using YOLO:

- A region of interest (ROI) is defined.
- The model is trained to identify this object.
- A cut is made in the region identified by YOLO to be used in the final segmentation step.

3.2.1. Defining ROI

To train YOLO, it is essential to define the ROI on the slices to be trained. However, the database only contains markings for the spinal cord. Therefore, to establish a region of interest, we use the markings made by spinal cord experts as a reference. Initially, we find the object's center of mass marked on the slice. Next, we draw a bounding box of dimensions 96×96 , with the same center of mass as the marking as its central pixel. We then plot the original slices from all training images to create the model in YOLO. Thus, from now on, we have a data set with the ROIs of the spine region. Furthermore, it is worth highlighting that these dimensions were chosen empirically and produced the best results.

3.2.2. YOLOv8 Training and Cut ROI

Based on the delimitation carried out in the previous step, the YOLOv8 model was used to detect and crop the region of interest automatically. This approach ensured that only the relevant area of the spine region was preserved, optimizing it for subsequent analyses without the need for prior notes or even the application of registration techniques (rigid or deformable) to the CT images.

The YOLOv8 model, developed by Ultralytics in 2015, is now being trained. Recognized for its effectiveness in object detection, this version of YOLO also stands out for its ability to support various computer vision tasks [Ultralytics 2024]. The choice of YOLOv8 for this activity was motivated by the need to detect and crop the ROI accurately. Its remarkable efficiency and accuracy in these tasks make it an ideal choice for this scenario. We chose the YOLOv8L version, a larger and more complex iteration of the model, due to its proven ability to handle more complex challenges and offer even greater accuracy.

Figure 3 illustrates the detection and cropping step with YOLOv8 [Ultralytics 2024].

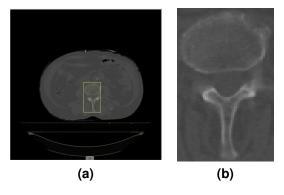


Figure 3. Spine region detection. (a) ROI detection, (b) ROI cropped.

Finally, the same cut is made to the corresponding images of the expert's markup to maintain the structure for the next step.

3.3. Spinal cord segmentation

After detecting and cropping the region of interest using YOLOv8, the resulting images are taken as input to the U-Net to segment the ME.

3.3.1. U-Net Training

U-Net is a convolutional neural network (CNN) architecture proposed by Olaf Ronneberger et al. in 2015, known for its good generalization in biomedical image segmentation. Its "U" shaped structure incorporates an encoder, which is responsible for reducing the dimensionality of the input, and a decoder, which increases the dimensionality to generate a detailed segmentation mask. The direct connection between the encoder and decoder layers preserves important contextual information during the downsampling and upsampling processes, contributing to the production of accurate segmentations [Ronneberger et al. 2015].

For training, the data is divided into training, validation, and testing. For training and validation, the ROI's recalled in the previous stage, both from the exam and the marking, were considered. During training, the loss function *dice-loss* was used as the aim was to optimize the segmentation closest to the training marking. For testing, the results were predicted and compared to the marking to evaluate the validation metrics (Sção 3.3.2). It is noteworthy that the U-Net architecture was trained according to the standard provided by the Keras library.

3.3.2. Evaluation of Results

In this section, after the U-Net prediction in section 3.3, metrics will be extracted to evaluate the performance of the proposed method. We used metrics commonly used in medical imaging problems, such as the Sørensen-Dice Index (Dice); Sensitivity, Specificity, and Accuracy; and the Area under the ROC Curve (AUC) [Diniz et al. 2021a].

4. Results and discussion

In this section, the results obtained by the proposed method for ME detection and segmentation will be presented. Furthermore, the results will be discussed, including a comparison with related work presented in section 2.

4.1. Training environment

To develop the method, the following hardware configurations were used: Processor: Ryzen 9 PRO 3900; RAM: 128GB DDR4; Video cards: 3x NVIDIA GeForce RTX 2080 TI 12GB.

The data was divided as provided by the challenge described in Section 3.1, where 7558 CTs were used for training and validation and 1975 CTs were used for testing. When training YOLO, the following hyperparameters were used: 200 epochs, two batch sizes, Stochastic Gradient Descent (SGD) as an optimizer, and a learning rate 0.01. It should be noted that these are the library's standard hyperparameters. In turn, for training the U-Net, the following hyperparameters were used: 200 epochs, eight batch sizes, Adam as an optimizer, and a learning rate of 0.0001.

4.2. Results for spine region detection

According to the dataset presented in Section 3.1, this section describes the results achieved by YOLO (Section 3.2.2.) for detecting and cropping the ROI. Below are the results of the metrics in YOLOv8 training:

Recall: 94.90%;Precision: 95.05%;mAP50: 96.60%.

Furthermore, Figure 4 shows the network training process, demonstrating significant accuracy in ROI detection. It displays important metrics such as precision, recall, and mAP50, as well as error rates during training and validation.

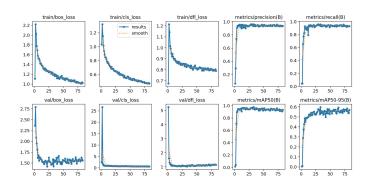


Figure 4. Performance metrics during YOLOv8 training.

The choice of YOLO proved robust for detecting the spine region, reducing the scope necessary for the next stage of ME segmentation.

4.3. Results for ME segmentation

With the column region delimited in Section 3.2.2, we proceed with the ME segmentation using U-Net in Section 3.3. To this end, we highlight some research points regarding the proposed method and solution. We will detail the conduct of the experiments in this stage, highlighting the following points:

- Some CT scan images, especially from the head to the spinal region, were observed without markings, as shown in Figure 5.
- Another crucial aspect observed is the lack of marking between slices, as shown in Figure 6. In other words, some slices are within a range that should have marking, but these are not marked. In this way, [Huang et al. 2020, Smith et al. 2019, Wu et al. 2018, LeCun and et al. 2021, Shorten and Khoshgoftaar 2019, Litjens and et al. 2017, Ravi and et al. 2020] provide valuable insights into the challenges faced in medical image segmentation, highlighting the importance of identifying and addressing these issues to improve analysis methods and image processing.

These points highlight the relevance of data preprocessing to address specific issues encountered in medical image segmentation tasks. Our method seeks to integrate effective solutions to improve the accuracy and reliability of ME segmentation.

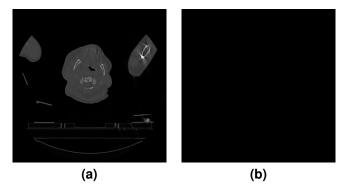


Figure 5. Head image. (a) Examination slice, (b) Marking slice without any ME delimitation.

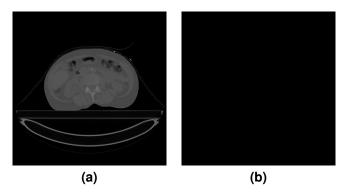


Figure 6. Column image. (a) Examination slice, (b) Marking slice without any ME delimitation.

Given this problem, we organized three experiments to validate the proposed method and improve ME segmentation:

- **Experiment I:** In this experiment, we performed training and prediction with U-Net on the entire dataset, including images, without expert marking.
- **Experiment II:** In the second experiment, we only removed images without expert markings from the training and validation set. After this removal, We retrained the network, keeping the test set with all images.
- Experiment III: Finally, in the third experiment, we removed the unmasked images from the test set using the network trained in experiment II.

Table 2 describes the results for each experiment.

Table 2. Results for experiments I, II and III

Experiment	Dice	Sensitivity	Specificity	Accuracy	AUC
I	76%	86%	98%	98%	92%
II	86%	86%	99%	98%	93%
III	88%	92%	99%	99%	96%

In Experiment I, where all images, including those without labels, were used, the Dice coefficient reached 76%. This suggests that unmarked slices introduced noise into the model training, resulting in less accurate segmentation of the ME. The absence of marking in some slices may have confused the model, leading to segmentation errors.

By removing the unlabeled images from the training and validation set in Experiment II, there was a significant improvement in the Dice coefficient, which increased to 86

In Experiment III, by removing unmarked and unmasked head images in the middle slices of the test set, there was a further performance improvement, with a Dice coefficient of 88%. This suggests that selectively removing these problematic images during the testing phase also played a crucial role in increasing segmentation accuracy. Testing the model on a clean dataset of problematic images made it possible to assess its ability to segment ME more accurately. It is worth noting that the fact that the model segments a region where there is no marking negatively influences the index. Therefore, this Experiment proves the method's robustness as it can be validated in marked slices.

In Figure 7, it is possible to compare the segmentation made by the expert, indicated in blue, with the neural network's prediction, indicated in red. In some cases, the segmentation performed by the expert and the neural network's prediction are quite similar, demonstrating the network's ability to produce results close to those of an expert. Furthermore, in situations with no expert marking, the neural network could segment the spinal cord accurately, demonstrating its ability to predict regions of interest even in the absence of previous markings. However, it is worth noting that despite consistent performance in many cases, there were situations where the neural network could not predict correctly due to gaps in marking. These results highlight both the robustness of the proposed method and the continued need to improve the model to deal with variations and gaps in marking.

In Figure 8, we present a case study that illustrates a slice of a patient's head region, as mentioned in Figure 5. In this section, it is important to note that the specialist

does not carry out any marking. The absence of this initial marking may suggest the specialist's difficulty when manually segmenting the region of interest. However, even in the face of this difficulty, the neural network could accurately predict the region of interest.

This method's ability to perform precise and complete segmentations, even in the absence of prior markings by the specialist, is remarkable. This is evidenced in Figure 7, where the neural network was able to segment the spinal cord in similar situations, filling gaps in the segmentation and proving robust in the absence of previous markings.

These examples highlight the proposed model's effectiveness and ability to perform reliable segmentations in various clinical scenarios. Furthermore, they demonstrate the method's usefulness as a valuable aid for medical specialists. The method's ability to overcome difficulties encountered in manual marking and its accuracy and robustness suggest that it could play a significant role in supporting medical diagnosis and treatment.

In summary, the presence of slices without expert marking significantly impacted the results, impairing the accuracy of ME segmentation. However, selectively removing these images during data preprocessing resulted in considerable improvements in model performance, highlighting the importance of careful input data management to achieve more accurate and reliable results.

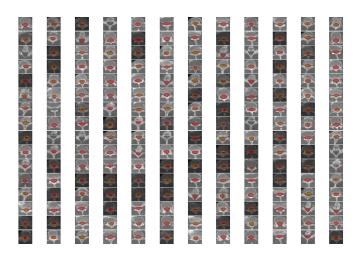


Figure 7. Case study: some slices of a patient from the testing dataset.

4.4. Discussion

By comparing the results of the proposed experiments with the results of related studies, it is possible to observe a general trend of similar or superior performance in the proposed methods. The Dice metric, in particular, is a crucial measure of overlap between segmented and true masks and is widely used in evaluating segmentation algorithms in medical images. The table below compares Dice values between the proposed experiments and related studies.

Comparing the results of the experiments with related works, a significant evolution in the performance of the organ segmentation method in computed tomography (CT) images is observed. Initially, the proposed method achieved a Dice of 76% in Experiment 1, which was below some of the values reported in the literature, such as the 78.20% in

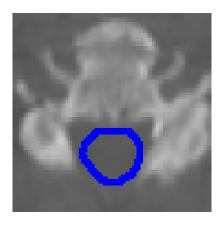


Figure 8. Case study: a slice of a patient from the head region that does not have the specialist's marking and the network predicted the region of interest.

Table 3. Comparison of Results with related works.

Method	Dice (%)
Diniz et al. [Diniz et al. 2021b]	78.20
Roncaglioni et al. [Roncaglioni 2021]	81.69
Lambert et al. [Lambert et al. 2020]	-
Cao et al. [Cao et al. 2021]	91
Chen et al. [Chen et al. 2021]	88
Spinoni et al. [SPINONI 2021]	84.05
Feng et al. [Feng et al. 2019]	89
Vieira et al. [Vieira et al. 2022]	86
Proposed - Experiment 1	76
Proposed - Experiment 2	86
Proposed - Experiment 3	88

the study by Diniz et al. and the 81.69% from Roncaglioni et al. However, by implementing improvements in later experiments, there was a notable progression. In Experiment 2, Dice increased to 86%, surpassing the results of some related studies. In Experiment 3, the method reached a Dice of 88%. The proposed method does not outperform Feng et al. and Cao et al.; however, these two works still need to address the problem of prior detection to eliminate unnecessary structures from the CT; in addition, they used 3D approaches that demand more computational resources.

Thus, the proposed method is believed to have the potential to become a valuable contribution to the field of medical image analysis and processing, offering significant benefits for clinical practice and the treatment of patients with cancer and other medical conditions.

5. Conclusion

Based on the results obtained in the experiments carried out, we observed a consistent evolution in several evaluation metrics, such as Dice, Sensitivity, and AUC. These improvements indicate a gradual optimization of the model, resulting in an improved ability to perform precise segmentation of the ME, which occupies a prominent space in the

literature.

However, the results also suggest that the proposed method can be improved further. For future work, an interesting approach would be to explore regularization or data augmentation techniques to improve model performance further. Furthermore, investigating the use of YOLOv8 to detect and segment the region of interest directly in an end-to-end approach could be a promising alternative, eliminating the need to use U-Net for segmentation.

References

- Cao, Z., Yu, B., Lei, B., Ying, H., Zhang, X., Chen, D. Z., and Wu, J. (2021). Cascaded se-resunet for segmentation of thoracic organs at risk. *Neurocomputing*, 453:357–368.
- Chen, X., Sun, S., Bai, N., Han, K., Liu, Q., Yao, S., Tang, H., Zhang, C., Lu, Z., Huang, Q., et al. (2021). A deep learning-based auto-segmentation system for organs-at-risk on whole-body computed tomography images for radiation therapy. *Radiotherapy and Oncology*, 160:175–184.
- Diniz, J., Ferreira, J., Silva, G., Quintanilha, D., Silva, A., and Paiva, A. (2021a). Segmentação de coração em tomografias computadorizadas utilizando atlas probabilístico e redes neurais convolucionais. In *Anais do XXI Simpósio Brasileiro de Computação Aplicada à Saúde*, pages 83–94, Porto Alegre, RS, Brasil. SBC.
- Diniz, J., Silva, A., and Paiva, A. (2021b). Methods for segmentation of spinal cord and esophagus in radiotherapy planning computed tomography. In *Anais Estendidos do XXXIV Conference on Graphics, Patterns and Images*, pages 21–27, Porto Alegre, RS, Brasil, SBC.
- Feng, X., Qing, K., Tustison, N. J., Meyer, C. H., and Chen, Q. (2019). Deep convolutional neural network for segmentation of thoracic organs-at-risk using cropped 3d images. *Medical physics*, 46(5):2169–2180.
- Huang, C., Boucneau, T., Theaud, G., Tourdias, T., and Rousseau, F. (2020). Automatic segmentation of the spinal cord and intramedullary multiple sclerosis lesions with convolutional neural networks. *NeuroImage*, 210:116555.
- Lambert, Z., Petitjean, C., Dubray, B., and Kuan, S. (2020). Segthor: Segmentation of thoracic organs at risk in ct images. In 2020 Tenth International Conference on Image Processing Theory, Tools and Applications (IPTA), pages 1–6. IEEE.
- Langner, U. W., Sohn, H., Hammon, M., Youn, I., Knechtges, P. M., Pahwa, S., Beesley, L., Peitz, S., Schachtschneider, P., Neumann, R., and Others (2013). The cancer imaging archive (tcia): A large-scale reference image database for imaging research. *Journal of Digital Imaging*, 26(3):613–618.
- LeCun, Y. and et al. (2021). Practical considerations in medical image segmentation: Deep learning-based approaches. *Nature Reviews Materials*, 6:205–218.
- Lima, L. and Oliveira, M. (2018). Using deep learning for classification of early lung nodules on computed tomography images. In *Anais do XVIII Simpósio Brasileiro de Computação Aplicada à Saúde*, Porto Alegre, RS, Brasil. SBC.
- Litjens, G. and et al. (2017). A survey on deep learning in medical image analysis. *Medical Image Analysis*, 42:60–88.

- Matos, C., Oliveira, M., Diniz, J., Fernandes, A., Junior, G. B., and Paiva, A. (2023). Ppm-deeplab: Módulo de pirâmide de pooling como codificador da rede deeplabv3+ para segmentação de rins, cistos e tumores renais. In *Anais do XXIII Simpósio Brasileiro de Computação Aplicada à Saúde*, pages 210–221, Porto Alegre, RS, Brasil. SBC.
- McKeown, S. R., Hatfield, P., Prestwich, R. J., Shaffer, R. E., and Taylor, R. E. (2015). Radiotherapy for benign disease; assessing the risk of radiation-induced cancer following exposure to intermediate dose radiation. *British Journal of Radiology*, 88(1056):20150405.
- Morbidity, T. (2016). Organs at risk and morbidity-related concepts and volumes. *J. ICRU*, 13(1–2):Rep. 89.
- Ravi, R. and et al. (2020). A review on deep learning techniques for brain tumor segmentation and classification. *IEEE Access*, 8:109114–109130.
- Roncaglioni, P. (2021). Ensemble methods for multi-organ semantic segmentation.
- Ronneberger, O., Fischer, P., and Brox, T. (2015). U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer.
- Santos, P., Brito, V., Filho, A. C., Sousa, A., Diniz, J., and Luz, D. (2023). Efficientbacillus: uma arquitetura profunda para detecção dos bacilos de koch. In *Anais do XXIII Simpósio Brasileiro de Computação Aplicada à Saúde*, pages 198–209, Porto Alegre, RS, Brasil. SBC.
- Shirato, H. and et al (2018). Selection of external beam radiotherapy approaches for precise and accurate cancer treatment. *Journal of Radiation Research*, 59(suppl_1):i2–i10.
- Shorten, S. and Khoshgoftaar, T. M. (2019). Data augmentation techniques for medical image analysis. *IEEE Access*, 7:111783–111796.
- Smith, J., Johnson, E., and Williams, R. (2019). Automated segmentation of medical images: a rapid method for segmenting kidney. *Medical Image Analysis*, 52:134–145.
- SPINONI, R. (2021). Multi organ semantic segmentation in ct scans.
- Ultralytics (Acessado em 2024). Yolov8: A Última evolução do yolo pela ultralytics. https://docs.ultralytics.com/#where-to-start.
- Vieira, P., Vogado, L., Lopes, L., Lira, R., Neto, P. S., Magalhaes, D., and Silva, R. (2022). Detecçao de doenças em imagens de raios-x da coluna lombo-sacra com convnets. In *Anais do XXII Simpósio Brasileiro de Computação Aplicada à Saúde*, pages 299–310. SBC.
- West, C. P., Dyrbye, L. N., and Shanafelt, T. D. (2018). Physician burnout: Contributors, consequences and solutions. *Journal of Internal Medicine*.
- Wu, J., Gong, G., Cao, L., Czito, B. G., Liu, Z., Qiu, X., and Yin, F.-F. (2018). Automated segmentation of thoracic and lumbar vertebrae on chest ct for image-guided radiation therapy of spine cancer. *Medical physics*, 45(12):5497–5506.