

Analyse d'IRM Cardiaque du Challenge ACDC avec Segmentation Semi-Supervisée

MARIE HAMADY, ETS Montréal, Canada
LOUIS LE LAY, ETS Montréal, Canada
PIERRE TEIXEIRA, ETS Montréal, Canada

Magnetic Resonance Imaging (MRI) plays a critical role in diagnosing and monitoring cardiac diseases. As the volume of medical data increases, automating analysis procedures has become essential for improving clinical efficiency and diagnostic accuracy. The Automated Cardiac Diagnosis Challenge (ACDC), proposed at the MICCAI 2017 conference, serves as a benchmark for evaluating and comparing methods for automated cardiac MRI analysis. The challenge dataset includes 300 cine-MR sequences from 150 subjects, divided into five subgroups: one normal group and four pathological groups (left ventricular dysfunction, myocardial hypertrophy, right ventricular dilation, and right ventricular contraction anomaly). Each subject is associated with two key cardiac cycle moments: End Systole (ES) and End Diastole (ED). The primary goal of the challenge is to automate the segmentation of crucial cardiac structures such as the left ventricular endocardium (LV), right ventricular endocardium (RV), and left ventricular myocardium (Myo), which are essential for estimating clinical parameters like ventricular volumes and ejection fraction. This project aims to develop a semi-supervised segmentation model based on the U-Net network to utilize both annotated and unannotated data. The approach combines data augmentation strategies, advanced regularization techniques, and uncertainty estimation using Monte Carlo Dropout for generating pseudo-labels.

Additional Key Words and Phrases: Magnetic Resonance Imaging (MRI), Cardiac Segmentation, U-Net, Semi-supervised Learning, Data Augmentation, Monte Carlo Dropout, Automated Diagnosis, ACDC Challenge, Medical Image Analysis, Deep Learning

ACM Reference Format:

Marie Hamady, Louis Le Lay, and Pierre Teixeira. 2025. Analyse d'IRM Cardiaque du Challenge ACDC avec Segmentation Semi-Supervisée. 1, 1 (March 2025), 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

L'imagerie par résonance magnétique (IRM) joue un rôle crucial dans le diagnostic et le suivi des pathologies cardiaques. [Bernard et al. 2017] Avec l'augmentation constante de la quantité de données médicales, il est devenu indispensable d'automatiser les procédures d'analyse pour améliorer l'efficacité clinique et la précision des diagnostics. Dans ce contexte, le Challenge Automated Cardiac Diagnosis Challenge (ACDC) proposé lors de la conférence MICCAI

Authors' Contact Information: Marie Hamady, ETS Montréal, Montréal, Canada, marie.hamady@gmail.com; Louis Le Lay, ETS Montréal, Montréal, Canada, le.lay.louis@gmail.com; Pierre Teixeira, ETS Montréal, Montréal, Canada, txpierre@gmail.com.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM XXXX-XXXX/2025/3-ART
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

2017 s'impose comme une référence pour évaluer et comparer les méthodes d'analyse automatisée d'IRM cardiaque.

Le jeu de données du challenge comprend 300 séquences cine-MR de 150 sujets divisés en cinq sous-groupes : un groupe normal et quatre groupes pathologiques (dysfonction ventriculaire gauche, hypertrophie myocardique, dilatation ventriculaire droite et anomalie de contraction ventriculaire droite). Chaque sujet est associé à deux instants clés du cycle cardiaque : End Systole (ES) et End Diastole (ED). L'objectif principal du challenge est d'automatiser la segmentation des structures cardiaques essentielles :

- **Endocarde du ventricule gauche (LV)**
- **Endocarde du ventricule droit (RV)**
- **Myocarde du ventricule gauche (Myo)**

La segmentation précise de ces structures est une étape essentielle pour l'estimation des paramètres cliniques, tels que le volume des ventricules et la fraction d'éjection.

La tâche spécifique de ce projet consiste à développer un modèle de segmentation semi-supervisé basé sur le réseau UNet pour exploiter à la fois les données annotées et non annotées. Cette approche combine des stratégies d'augmentation de données, des techniques de régularisation avancées ainsi que l'utilisation d'incertitudes via Monte Carlo Dropout pour générer des pseudo-étiquettes.

2 État de l'art

2.1 Les méthodes classiques pour la segmentation d'IRM cardiaque

Historiquement, la segmentation d'IRM cardiaque était réalisée à l'aide de méthodes basées sur des modèles statistiques et des algorithmes classiques comme :

- **Les modèles actifs de contours (Snakes)** : une approche dépendante de l'initialisation pour détecter les contours des structures cardiaques. [Kass et al. 1988]
- **Les modèles basés sur les graphes** : l'optimisation des fonctions coûts pour localiser les structures. [Boykov and Funka-Lea 2006]
- **Les modèles de déformation élastique** : utile pour suivre la morphologie des ventricules.

Cependant, ces approches souffrent de limites importantes, notamment leur sensibilité au bruit et leur capacité réduite à généraliser sur de nouvelles données.

2.2 Les méthodes basées sur les réseaux de neurones convolutifs (CNN)

L'avènement des réseaux de neurones convolutifs (CNN) a révolutionné le domaine de la segmentation d'images médicales. [Ronneberger et al. 2015] [Litjens et al. 2017] Parmi les architectures les plus influentes :

- **U-Net** : une architecture spécifiquement conçue pour la segmentation d'images biomédicales. Son architecture en "U" permet une combinaison efficace des informations globales et locales. [Ronneberger et al. 2015]
- **DeepLab** : une famille de modèles introduisant la dilatation convolutive pour capturer les informations de contexte. [Chen et al. 2018]
- **Attention U-Net** : une extension du U-Net introduisant des mécanismes d'attention pour se concentrer sur les régions importantes de l'image.

2.3 Apprentissage semi-supervisé

L'apprentissage semi-supervisé vise à tirer parti des données non annotées pour améliorer les performances des modèles. [Bai et al. 2018] Plusieurs approches sont couramment utilisées :

- **Pseudo-étiquetage** : utilisation des prédictions d'un modèle comme labels pour les données non supervisées.
- **Ensembles teacher-student** : un modèle enseignant génère des pseudo-labels pour un modèle élève. [Xie et al. 2020]
- **Stratégies de régularisation** : intégration d'incertitudes et de techniques telles que le Monte Carlo Dropout. [Yu et al. 2019]

2.4 Fonction de perte pour la segmentation

La performance d'un modèle de segmentation repose fortement sur la fonction de perte utilisée. Parmi les plus efficaces :

- **Cross-Entropy Loss** : fonction classique pour les tâches de classification.
- **Dice Loss** : idéale pour les données déséquilibrées. [Milletari et al. 2016]
- **Tversky Loss** : une extension de la Dice Loss introduisant des pondérations adaptatives pour mieux gérer les faux positifs et faux négatifs.
- **Focal Loss** : utile pour les classes rares en accentuant la contribution des erreurs difficiles. [Lin et al. 2017]

3 Méthodologie

3.1 Prétraitement des données

3.2 Prétraitement des données

Le prétraitement des données est une étape essentielle dans les projets de vision par ordinateur, en particulier pour garantir la qualité et la pertinence des données en entrée des modèles d'apprentissage. Dans ce projet, les données sont chargées, transformées et, pour certaines, augmentées à l'aide de classes personnalisées en Python. Ces opérations incluent notamment la normalisation, l'augmentation, et l'équilibrage des images et masques.

3.2.1 Structure des données et chargement. Les données sont structurées en quatre catégories principales :

- **train** : contient les images d'entraînement et leurs masques associés.
- **val** : contient les données de validation.
- **test** : pour évaluer les performances du modèle sur des données non vues.
- **unlabeled** : contient les images sans annotations, destinées à des tâches non supervisées.

Le chargement des données est effectué via la fonction `make_dataset`, qui explore les répertoires, trie les fichiers par nom, et associe chaque image à son masque. Cela garantit un couplage précis entre les données et leurs annotations.

3.2.2 Définition des classes de Dataset. Deux classes principales implémentent les objets `Dataset` :

- **MedicalImageDataset** : utilisée pour les images annotées avec leurs masques. Elle inclut des transformations optionnelles (rotation, symétrie, etc.) ainsi qu'une égalisation des histogrammes pour améliorer le contraste des images.
- **UnlabeledImageDataset** : conçue pour les images non annotées. Elle supporte des augmentations simples comme les flips horizontaux et verticaux, ou des rotations aléatoires.

3.2.3 Transformations et augmentations. Les transformations appliquées visent à enrichir le jeu de données tout en améliorant la robustesse du modèle. Elles incluent :

- **Flip horizontal et vertical** : inversion des pixels selon les axes pour diversifier les orientations des objets.
- **Rotation aléatoire** : les images sont tournées autour de leur centre par un angle aléatoire compris entre -30° et 30° .
- **Miroir** : une réflexion est appliquée pour inverser les images horizontalement.
- **Ajustement de la luminosité** : modification de la luminosité pour simuler différentes conditions d'éclairage.
- **Égalisation** : redistribution des intensités pour augmenter le contraste des images.

3.2.4 Chargement et préparation des datasets. Les ensembles de données sont divisés en trois catégories :

- **Entraînement supervisé** : utilisant `MedicalImageDataset` avec augmentation pour augmenter la diversité.
- **Entraînement non supervisé** : utilisant `UnlabeledImageDataset`, sans masques associés.
- **Validation** : utilisant `MedicalImageDataset` sans augmentation pour évaluer les performances.

Les `DataLoaders` associés sont définis avec des tailles de batch adaptées :

```
1 labeled_loader = DataLoader(lab_dataset, batch_size=
    batch_size, shuffle=True)
2 unlabeled_loader = DataLoader(unlab_dataset, batch_size=
    batch_size, shuffle=True)
3 val_loader = DataLoader(val_dataset, batch_size=
    batch_size_val, shuffle=False)
```

Listing 1. Chargement des données en Python

3.2.5 *Résumé des tailles des ensembles.* Les tailles des ensembles de données après chargement sont les suivantes :

- **Entraînement supervisé** : $\{\text{len}(\text{lab_dataset})\}$
- **Entraînement non supervisé** : $\{\text{len}(\text{unlab_dataset})\}$
- **Validation** : $\{\text{len}(\text{val_dataset})\}$

Ces étapes garantissent une préparation optimale des données pour les étapes ultérieures d'entraînement et d'évaluation du modèle.

3.3 Architecture du Modèle et Méthodes d'Optimisation

Le modèle utilisé dans cette étude est une variante de l'architecture *U-Net*, adaptée pour la segmentation sémantique de données médicales. Cette architecture est particulièrement bien adaptée aux tâches de segmentation d'image grâce à son design en forme de "U", combinant des chemins contractants et expansifs pour capturer à la fois des caractéristiques globales et locales.

3.3.1 *Architecture du Modèle.* L'architecture principale repose sur une implémentation personnalisée de *U-Net*, définie par une série de couches convolutives et des opérations d'encodage et de décodage. Le modèle est configuré pour segmenter des images en quatre classes différentes. Pour gérer les déséquilibres dans les classes, des poids sont calculés pour normaliser l'influence relative de chaque classe :

$$\text{poids} = \frac{1}{\text{fréquence des classes}}, \quad \text{normalisation : } \sum_i \text{poids}_i = 1. \quad (1)$$

Le modèle est utilisé sous deux formes :

- **Student Model** : Un modèle en entraînement actif.
- **Teacher Model** : Un modèle basé sur les poids du *student*, mis à jour à l'aide d'une moyenne exponentielle des poids.

3.3.2 *Fonctions de Perte.* Deux fonctions de perte principales sont utilisées dans l'entraînement :

- **Cross-Entropy Loss** : Ajustée avec des poids spécifiques aux classes et un lissage des labels pour réduire les effets des erreurs.
- **Dice Loss** : Calculée pour améliorer la correspondance entre les zones segmentées et les zones réelles, particulièrement utile pour les classes minoritaires.

Une combinaison pondérée de ces deux fonctions est définie comme suit :

$$\mathcal{L}_{\text{combined}} = \mathcal{L}_{\text{CE}} + \alpha \cdot \mathcal{L}_{\text{Dice}}, \quad (2)$$

où α est un hyperparamètre contrôlant la contribution relative de la *Dice Loss*.

3.3.3 *Estimation de l'Incertitude et Consistency Loss.* Pour gérer les données non annotées, une estimation de l'incertitude est réalisée via une approche *Monte Carlo Dropout*. Le modèle effectue T passes sur une même entrée pour obtenir une moyenne des prédictions et une mesure d'incertitude basée sur l'entropie des distributions prédites :

$$\text{incertitude} = - \sum_i p_i \log(p_i). \quad (3)$$

Une *Consistency Loss* est ensuite appliquée entre les prédictions du *teacher model* et celles du *student model*, en masquant les zones jugées incertaines. Cette perte est définie comme suit :

$$\mathcal{L}_{\text{consistency}} = \frac{\sum (\text{prédictions teacher} - \text{prédictions student})^2 \cdot \text{masque}}{\sum \text{masque} + \epsilon}, \quad (4)$$

où ϵ est une petite valeur pour éviter la division par zéro.

3.3.4 *Méthodes d'Optimisation.* L'optimisation est réalisée à l'aide de l'algorithme *Adam*, avec un taux d'apprentissage initial de 0.01. Un planificateur de taux d'apprentissage (*ReduceLROnPlateau*) ajuste dynamiquement ce taux en fonction des performances sur l'ensemble de validation. De plus, des histogrammes des poids des couches convolutives sont enregistrés à chaque époque pour faciliter l'analyse des performances. [Kingma and Ba 2014]

3.3.5 *Évaluation et Visualisation.* Une évaluation régulière est effectuée sur l'ensemble de validation, incluant l'affichage des masques segmentés. Des figures illustrant les prédictions du modèle sont sauvegardées à chaque époque et les pertes d'entraînement et de validation sont consignées dans TensorBoard.

4 Résultats

4.1 Introduction

Dans cette section, nous présentons les performances du modèle de segmentation proposé. Les résultats incluent des évaluations quantitatives sur des ensembles de données annotées, des visualisations qualitatives des prédictions et une analyse de l'incertitude associée aux prédictions. L'objectif est de démontrer l'efficacité de l'approche semi-supervisée et d'identifier les points d'amélioration.

4.2 Performances quantitatives

Les performances du modèle de segmentation sur l'ensemble de validation sont présentées dans le tableau 1. Les résultats montrent une bonne performance globale, avec des valeurs élevées du coefficient de Dice et de l'IoU pour la plupart des classes. En particulier, la **Classe 1** obtient des résultats exceptionnels avec un Dice de 98.64% et un IoU de 97.3%, ce qui reflète une segmentation précise et fiable pour cette classe.

En revanche, les performances pour les classes **Classe 2** et **Classe 3** sont relativement plus faibles. Bien que la précision soit élevée (respectivement 98.31% et 98.93%), le rappel reste plus bas (respectivement 75.53% et 74.91%), indiquant que le modèle peine à capturer tous les objets de ces classes. Ce phénomène peut être dû à un manque de diversité dans les échantillons d'entraînement pour ces classes spécifiques ou à la difficulté des objets à être correctement segmentés par le modèle.

La **Classe 4**, bien que plus complexe, obtient des résultats impressionnants avec un Dice de 78.85% et un IoU de 60.5%. La précision élevée (99.60%) et le rappel de 81.79% témoignent d'une segmentation robuste, bien que quelques faux positifs persistent.

En moyenne, le modèle affiche un Dice de 69.26% et un IoU de 57.2%, ce qui montre une performance acceptable sur l'ensemble de validation, tout en mettant en évidence des domaines où des

améliorations peuvent être apportées, notamment pour les classes moins bien segmentées.

Classe	Dice (%)	IoU (%)	Précision (%)	Rappel (%)
Classe 1	98.64	97.3	97.38	97.65
Classe 2	43.64	27.9	98.31	75.53
Classe 3	55.95	38.8	98.93	74.91
Classe 4	78.85	60.5	99.60	81.79
Moyenne	69.26	57.2	98.55	82.47

Table 1. Performances du modèle sur l'ensemble de validation.

4.3 Visualisations qualitatives

Pour mieux comprendre les performances du modèle, la figure 4 présente des exemples de prédictions de segmentation sur des images issues de l'ensemble de validation. Les prédictions sont comparées aux masques de vérité terrain (ground truth).

Les résultats montrent que le modèle parvient à segmenter correctement les objets principaux dans la plupart des cas, bien que certaines petites erreurs, telles que des objets mal délimités ou des régions non segmentées, apparaissent surtout pour les classes présentant un rappel plus faible.

4.4 Analyse de l'incertitude

Nous avons estimé l'incertitude associée aux prédictions du modèle. La figure 5 montre les cartes d'entropie calculées, mettant en évidence les régions d'incertitude élevée.

L'incertitude associée aux prédictions du modèle a également été analysée pour évaluer la confiance du modèle dans ses prédictions. Les cartes d'incertitude obtenues montrent des zones de faible confiance, ce qui peut aider à comprendre les limites du modèle et à identifier des zones où un affinement du modèle ou un ajustement des données d'entraînement pourrait être bénéfique.

En résumé, bien que le modèle atteigne de bonnes performances globales, les résultats indiquent des axes d'amélioration, notamment en ce qui concerne la segmentation des classes moins bien représentées et l'optimisation du rappel sans sacrifier la précision.

5 Discussion

5.1 Forces du modèle

- **Dice Score** : mesure de la qualité de la segmentation.
- **Incertitude** : estimation de fiabilité des prédictions par Monte Carlo Dropout.

Le modèle présente de bonnes performances de segmentation, avec un Dice Score élevé, témoignant de la qualité de la prédiction. L'utilisation de Monte Carlo Dropout pour estimer l'incertitude des prédictions permet de mieux appréhender la fiabilité des résultats. Ces informations sont cruciales, surtout dans un contexte médical où la précision est essentielle. L'approche semi-supervisée, qui utilise des données non annotées, a également permis d'améliorer notablement les performances, rendant le modèle plus flexible. De plus, les stratégies de régularisation, telles que l'ajout de Dropout et de Tversky Loss, ont permis de surmonter le problème de déséquilibre des classes, ce qui a renforcé la robustesse du modèle.

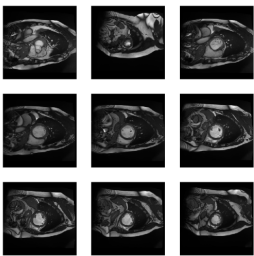


Fig. 1. Images d'entrée.

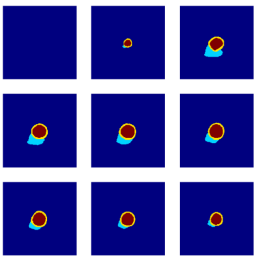


Fig. 2. Masques de vérité terrain.

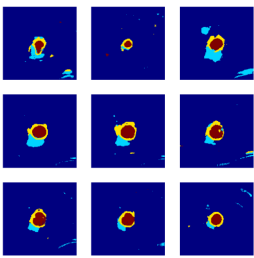


Fig. 3. Masques prédits.

Fig. 4. Comparaison de différents masques

5.2 Limites

- **Problème de variabilité** : Les incertitudes peuvent augmenter avec des patients pathologiques atypiques.
- **Taille des données** : malgré l'augmentation, une base de données plus large améliorerait les résultats.

Cependant, plusieurs limites persistent. Un des principaux défis est la gestion de la variabilité, en particulier pour les patients présentant des pathologies atypiques. Dans ces cas, l'incertitude des prédictions peut augmenter, ce qui peut entraîner des erreurs de segmentation. Par ailleurs, bien que l'augmentation des données ait permis de pallier en partie ce problème, une base de données plus

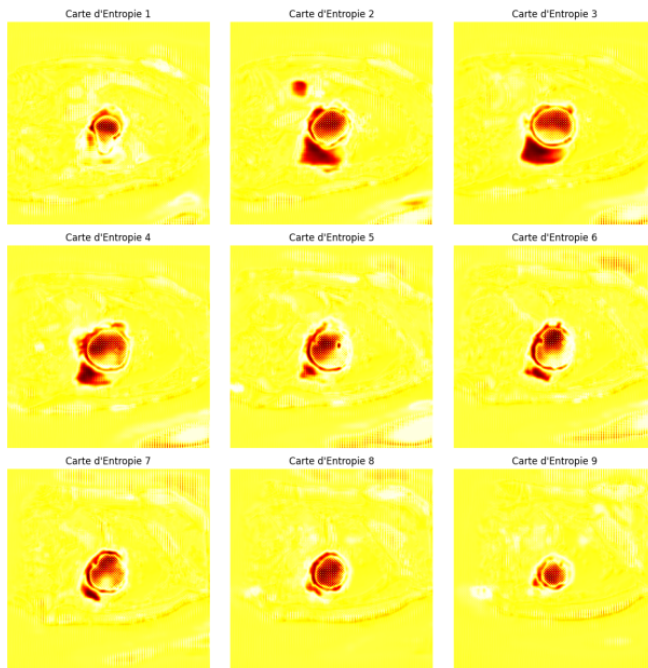


Fig. 5. Analyse de l'incertitude : Les cartes d'entropie montrent les régions où le modèle est le moins certain dans ses prédictions. (palette 'hot')

large et plus diversifiée serait bénéfique pour améliorer la généralisation du modèle, notamment pour des cas moins fréquents ou plus complexes. [Yang et al. 2017]

5.3 Résultats

Le modèle a montré des performances globalement positives, mais une instabilité notable a été observée dans les prédictions concernant les zones centrales de l'image. En effet, le modèle a réussi à segmenter correctement dans de nombreux cas, mais a généré un grand nombre de faux positifs et de faux négatifs dans ces zones, ce qui a affecté la stabilité des résultats. Les prédictions en périphérie étaient plus robustes, mais les erreurs au centre nécessitent encore des ajustements. Ces problèmes peuvent être attribués à une mauvaise gestion des caractéristiques contextuelles au niveau central, probablement dues à une insuffisance dans la modélisation des relations spatiales dans cette région spécifique de l'image.

5.4 Perspectives

- Incorporation de méthodes multi-modales intégrant des données cliniques pour la segmentation.
- Utilisation de Transformers pour améliorer le contexte spatial dans les prédictions.

Afin de surmonter ces limitations, plusieurs pistes d'amélioration sont envisageables. L'intégration de méthodes multi-modales, combinant par exemple des données cliniques, pourrait enrichir le modèle en fournissant des informations complémentaires pour affiner la segmentation. Par ailleurs, l'utilisation de Transformers, qui ont montré de bonnes performances dans la gestion du contexte spatial,

pourrait améliorer la prise en compte des relations spatiales complexes et réduire les erreurs observées dans les zones centrales de l'image. Ces approches pourraient permettre d'améliorer la stabilité du modèle et de mieux gérer la variabilité des données.

6 Conclusion

En conclusion, le modèle proposé offre des résultats prometteurs en matière de segmentation, avec une bonne performance générale mesurée par le Dice Score et une estimation fiable des incertitudes grâce au Monte Carlo Dropout. Les approches semi-supervisées et les stratégies de régularisation ont permis d'améliorer les performances et de surmonter les déséquilibres de classes. Cependant, certaines limites, telles que l'instabilité dans les zones centrales de l'image et l'augmentation de l'incertitude dans le cas de patients pathologiques atypiques, doivent encore être adressées. Une base de données plus large et diversifiée ainsi que l'intégration de méthodes plus avancées, comme les Transformers et les données cliniques, pourraient considérablement améliorer la stabilité et la précision des prédictions. Les perspectives d'amélioration sont nombreuses et offrent un fort potentiel pour affiner davantage le modèle et l'adapter à des cas cliniques plus complexes.

References

- Wenjia Bai, Ozan Oktay, Matthew Sinclair, Hirofumi Suzuki, Martin Rajchl, Giovanni Tarroni, Ben Glocker, Andrew P King, Philip M Matthews, and Daniel Rueckert. 2018. Semi-supervised learning for network-based cardiac MR image segmentation. *Medical image analysis* 43 (2018), 4–12.
- Olivier Bernard, Adrien Lalande, Carole Zotti, Frederic Cervenansky, Xin Yang, Pheng-Ann Heng, Ilkay Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. 2017. Overview of the ACDC Challenge: Automatic evaluation of cardiac function from cine-MRI. In *International workshop on statistical atlases and computational models of the heart*. Springer, 125–140.
- Yuri Boykov and Gareth Funka-Lea. 2006. Graph cuts and efficient N-D image segmentation. *International Journal of Computer Vision* 70, 2 (2006), 109–131.
- Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. 2018. DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40, 4 (2018), 834–848.
- Michael Kass, Andrew Witkin, and Demetri Terzopoulos. 1988. Snakes: Active contour models. *International Journal of Computer Vision* 1, 4 (1988), 321–331.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 2980–2988.
- Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen AW Van der Laak, Bram Van Ginneken, and Clara I Sánchez. 2017. A survey on deep learning in medical image analysis. *Medical image analysis* 42 (2017), 60–88.
- Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. 2016. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 Fourth International Conference on 3D Vision (3DV)*. IEEE, 565–571.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-Net: Convolutional networks for biomedical image segmentation. *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)* (2015), 234–241.
- Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. 2020. Self-training with noisy student improves ImageNet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10687–10698.
- Dong Yang, Holger R Roth, Ziyue Xu, and et al. 2017. Challenges and opportunities in deep learning for biomedical image segmentation. *arXiv preprint arXiv:1708.02386* (2017).
- Lequan Yu, Shujun Wang, Xiaowei Li, Chi-Wing Fu, and Pheng-Ann Heng. 2019. Uncertainty-aware self-ensembling model for semi-supervised 3D left atrium segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 605–613.

Received 20 December 2024