

RECIST WEAKLY SUPERVISED LESION SEGMENTATION VIA LABEL-SPACE CO-TRAINING

Lianyu Zhou¹, Dong Wei², Donghuan Lu², Wei Xue³, Liansheng Wang^{1*}, Yefeng Zheng²

¹Xiamen University; ²Tencent Jarvis Lab; ³OPPO

ABSTRACT

As an essential indicator for cancer progression and treatment response, tumor size is often measured following the response evaluation criteria in solid tumors (RECIST) guideline in CT slices. By marking each lesion with its longest axis and the longest perpendicular one, laborious pixel-wise manual annotation can be avoided. However, such a coarse substitute cannot provide a rich and accurate base to allow versatile quantitative analysis of lesions. To this end, we propose a novel weakly supervised framework to exploit the existing rich RECIST annotations for pixel-wise lesion segmentation. Specifically, a pair of under- and over-segmenting masks are constructed for each lesion based on its RECIST annotation and served as the label for co-training a pair of subnets, respectively, along with the proposed label-space perturbation induced consistency loss to bridge the gap between the two subnets and enable effective co-training. Extensive experiments are conducted on a public dataset to demonstrate the superiority of the proposed framework regarding the RECIST-based weakly supervised segmentation task and its universal applicability to various backbone networks.

Index Terms— Response evaluation criteria in solid tumors (RECIST), weakly supervised segmentation, label-space co-training

1. INTRODUCTION

Tumor size measurement in medical imaging and follow-ups is a widely accepted protocol for cancer monitoring [1]. In current clinical routine, the measurement is often performed in computed tomography (CT) slices manually by trained specialists, following the response evaluation criteria in solid tumors (RECIST) guideline [2]. Specifically, a couple of RECIST diameters are marked for each lesion (Fig. 1(b)), with the major diameter measuring the longest axis of the lesion and the minor measuring the longest perpendicular axis to the major, in the axial slice where the lesion appears largest. The RECIST diameters are commonly adopted as a time-efficient alternative to full lesion segmentation (e.g., stroking along the lesion boundary precisely or pixel-wise annotation), i.e., a

trade-off between measurement accuracy and annotation effort. As a coarse substitute, however, they cannot provide as much information as the full segmentation, which would allow versatile quantitative analysis such as area estimation or morphological properties.

Despite its great value, manual segmentation of medical images is notoriously laborious, tedious, error-prone, and subject to inter-observer variability. Therefore, tremendous efforts have been made to develop automated methods [3]. Significant breakthroughs have been achieved in the past decade by deep convolutional neural networks (DCNNs) [4]. However, full supervision of the DCNNs requires pixel-wise annotations of large datasets [5], which can be difficult or costly to obtain in practice. To reduce this burden, weak supervision in the form of simplified or partial labels, e.g., bounding boxes [6] or scribbles [7], has been studied with great interests. For lesion segmentation, the existing RECIST annotations are a natural and rich source of weak supervision.

Recently, several works have emerged that utilized the RECIST annotations for weakly supervised lesion segmentation [8, 9, 10, 11, 12, 13]. Despite the remarkable progress, we identify three common drawbacks of existing works that should be overcome to push the research towards better clinical applicability. First, almost all the existing methods [8, 9, 10, 11] relied on the GrabCut algorithm [14]—a classical unsupervised segmentation algorithm—to generate initial pseudo ground truth from the RECIST diameters for training the segmentation networks. As the quality of the pseudo ground truth largely depended on the initial seeds for GrabCut, the seeding strategy must be carefully devised to optimize the pseudo ground truth, which may encounter difficulty in generalization in practice; besides, the delicate pre-generation stage added unnecessary complexity to the training process, and even so, these empirical strategies were not guaranteed to always generate lesion masks with high fidelity. Second, many existing works involved iterative training procedures, where the pseudo ground truth was updated and the network performance gradually improved in rounds [9, 10], leading to exceedingly time-consuming training process. Third, most previous works assumed that the lesion-of-interest (LOI; an enlarged bounding box of the lesion) region was cropped out beforehand as input [8, 9, 10, 13]. Such assumption, while acceptable when developing early-stage

*Correspondence: lswang@xmu.edu.cn

L. Zhou, D. Wei and D. Lu contributed equally.

prototypes, may pose an obstacle to practical application as the clinicians still need to obtain the LOI first. A straightforward solution is to prefix a lesion detection model to the segmentation network [11]. However, such a two-stage structure can be unnecessarily redundant if a simpler one-stage model can achieve the same performance.

In this work, we present a novel one-stage co-training framework for RECIST supervised lesion segmentation in CT slices, which effectively addresses all the three drawbacks. Rather than over-tuning a classical unsupervised method for optimal pseudo ground truth, we instead make an intuitive observation that the RECIST diameters naturally compose two sets of masks: a quadrilateral connecting, and a circle circumscribing the four end points. Thus, we train a model whose two subnets are supervised with either of the two masks, respectively, Noting that the two masks are inherently under- and over-segmented representations of the lesion, respectively, we obtain the final prediction by averaging the corresponding under- and over-segmenting predictions. In this way, the training label construction is simple and straightforward. Inspired by the recent progress in self supervision with contrastive learning [15], we propose a novel consistency loss that contrasts the two subnets’ predictions for label-space perturbation based co-training [16]. Owing to the simplicity, robustness, and efficacy of the proposed dual label construction and label-space co-training, our framework can accept whole CT slices as input and get rid of the iterative refinement, while still being able to produce high-quality lesion segmentation. Last but not least, our framework is model-agnostic and can be readily applied to any standard segmentation backbone. Experimental results on a public dataset demonstrate the advantages of our framework over existing approaches.

2. METHOD

Problem Formulation. Following the literature, the target of this work is to train a model with RECIST annotations to perform accurate dense pixel-wise classification of lesion and non-lesion on axial CT slices. Formally, we view a slice and its lesion mask as two K -dimensional vectors: $\mathbf{I} \in [0, 1]^K$ and $\mathbf{M} \in \{0, 1\}^K$, respectively, where all pixels of a slice constitute an index set $\mathcal{P} = \{p | p = 0, 1, \dots, K - 1\}$. For a pixel p , $M_p = 1$ indicates that it is foreground (a lesion pixel), whereas $M_p = 0$ indicates background; all foreground pixels constitute a set $\mathcal{M} = \{p | M_p = 1, p \in \mathcal{P}\}$. Hence, \mathbf{M} and \mathcal{M} are alternative representations of the same mask. Similarly, the RECIST annotation of a lesion can be represented as a mask \mathbf{R} and a corresponding index set \mathcal{R} , where $R_p = 1$ indicates the pixel p is on the RECIST diameters and 0 otherwise. Therefore, given a training slice set with RECIST annotations¹ $D^{\text{train}} = \{(\mathbf{I}, \{\mathbf{R}\})\}$, the target is to train a segmentation model that can predict accurate $\hat{\mathbf{M}}$ for any

¹A slice may contain multiple lesions and thus has a set of annotations.

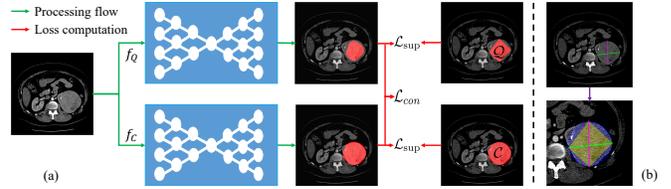


Fig. 1. (a) Overview of the proposed framework. (b) Illustration of the dual mask construction. Top: a CT slice with RECIST diameters overlaid. Bottom: the same slice zoomed for better visibility, where yellow color indicates definite foreground (the quadrilateral \mathcal{Q}), blue indicates the uncertain area \mathcal{A} , their combination indicates the circumscribed circle \mathcal{C} , and the rest color-less region indicates definite background.

unannotated test slice.

Dual Mask Construction from RECIST Diameters. In this work, we propose straightforward construction of a pair of simple yet contrasting masks for network supervision (Fig. 1(a)): 1) a quadrilateral that connects the four endpoints of the RECIST diameters with enclosing pixels forming an index set \mathcal{Q} for mask \mathbf{Q} : $\mathcal{Q} = \{p | Q_p = 1, p \in \mathcal{P}\}$, and 2) the minimum circumscribed circle of the diameters with enclosing pixels forming an index set \mathcal{C} for mask \mathbf{C} : $\mathcal{C} = \{p | C_p = 1, p \in \mathcal{P}\}$. Since most RECIST-measurable lesions show convex shapes [2], it is reasonable to analyze the general properties of \mathbf{Q} and \mathbf{C} assuming lesions with convex outlines [17]. Based on this assumption, it is easy to see that the quadrilateral is completely contained in the lesion, i.e., $\mathcal{Q} \subseteq \mathcal{M}$. Therefore, \mathbf{Q} is an under-segmentation of the genuine lesion mask \mathbf{M} . On the contrary, \mathbf{C} is an over-segmentation of \mathbf{M} , i.e., $\mathcal{M} \subseteq \mathcal{C}$. In fact, this circle is also the minimum circumscribed circle of the lesion, i.e., the smallest circle that can fully contain the lesion while including as little background as possible, given the lesion is convex. This is favorable especially compared to the bounding boxes of the RECIST diameters [18], which guarantees neither full inclusion of the lesion nor least inclusion of background. With the obvious “flaws” of being under- and over-segmenting, directly training with either of these two seemingly naive masks would lead to suboptimal results. Next, we describe how to effectively utilize them to build our framework.

Co-Training with Label-Space Perturbation Induced Consistency Loss. As show in Fig. 1(a), our framework mainly consists of two subnetworks f_Q and f_C supervised by \mathbf{Q} and \mathbf{C} , respectively, with its loss function defined as:

$$\mathcal{L}_{\text{sup}} = \ell_{\mathcal{P}}(\hat{\mathbf{Q}}, \mathbf{Q}) + \ell_{\mathcal{P}}(\hat{\mathbf{C}}, \mathbf{C}), \quad (1)$$

where $\hat{\mathbf{Q}} = f_Q(\mathbf{I})$ and $\hat{\mathbf{C}} = f_C(\mathbf{I})$ are the predicted probability maps by the two networks for a slice \mathbf{I} , respectively, and $\ell_{\mathcal{P}}(\cdot, \cdot)$ is a segmentation loss (such as the Dice loss) effective in set \mathcal{P} . \mathcal{L}_{sup} drives f_Q and f_C to output predictions mimicking their supervision masks \mathbf{Q} and \mathbf{C} , respectively. Therefore, the predictions are expected to be under-

and over-segmenting. A natural approach to fusing them is an averaging ensemble: $\bar{M} = (\hat{Q} + \hat{C})/2$. As later shown in our experiment, \bar{M} is consistently better than both \hat{Q} and \hat{C} . This is reasonable: as \hat{Q} and \hat{C} are biased estimations of the underlying ground truth towards diverging directions, averaging them is supposed to effectively cancel the biases. However, such a setting is essentially a simple ensemble of two networks trained synchronously but independently. To bridge the two networks, we propose to co-train f_Q and f_C by explicitly enforcing the consistency between \hat{Q} and \hat{C} with a contrasting loss:

$$\mathcal{L}_{\text{con}} = \ell_{\mathcal{P}}(\hat{Q}, \hat{C}). \quad (2)$$

The rationale is that Q and C are in fact intentionally introduced perturbations to the same underlying ground truth masks. Enforcing the consistency in the two subnets’ predictions across the perturbations can make them gradually approach a proper balance between the under- and over-segmenting biases, thus producing a more accurate segmentation. With a weight factor λ , combining \mathcal{L}_{sup} with \mathcal{L}_{con} yields the complete optimization target of our framework:

$$\mathcal{L} = \ell_{\mathcal{P}}(\hat{Q}, Q) + \ell_{\mathcal{P}}(\hat{C}, C) + \lambda \ell_{\mathcal{P}}(\hat{Q}, \hat{C}). \quad (3)$$

Region-Constrained Consistency Loss. Referring to Fig. 1(b), it is easy to derive that: $\forall p \in \mathcal{Q}, Q_p = M_p = C_p = 1$, and $\forall p \notin \mathcal{C}, Q_p = M_p = C_p = 0$, thus we have

$$\forall p \in \mathcal{Q} \cup (\mathcal{P} - \mathcal{C}), Q_p = M_p = C_p. \quad (4)$$

This is to say, both Q and C are consistent with the underlying ground truth lesion mask M in regions inside \mathcal{Q} (definite foreground) and outside \mathcal{C} (definite background). In our design, \mathcal{L}_{con} is targeted at the ambiguous region where Q and C disagree and the real ground truth is unknown, which is the region outside \mathcal{Q} but inside \mathcal{C} : $\mathcal{A} = \mathcal{C} - \mathcal{Q}$. On one hand, when f_Q and f_C are trained well by Q and C , respectively, their predictions should be the same in the agreed region $\mathcal{Q} \cup (\mathcal{P} - \mathcal{C})$, and \mathcal{L}_{con} does not function. On the other hand, if either one of them makes a mistake in this region, the mistake may mislead the other network through \mathcal{L}_{con} . Therefore, we further propose to improve \mathcal{L}_{con} by constraining its effective region within \mathcal{A} instead of the whole slice \mathcal{P} , thus Eq. (3) becomes:

$$\mathcal{L} = \ell_{\mathcal{P}}(\hat{Q}, Q) + \ell_{\mathcal{P}}(\hat{C}, C) + \lambda \ell_{\mathcal{A}}(\hat{Q}, \hat{C}), \quad (5)$$

which is expected to be more efficient, effective, and robust.

Relation with Related Work. Conventional co-training often co-trained the subnets with distinct inputs (e.g., different views of a web page [16] or MRI sequences in an exam [19]). In this work, in an attempt to utilize the inherent uncertainty of the RECIST annotation, we instead co-train the two subnets with two distinct supervision masks which are purposely constructed with diverging biases. This is in contrast with the GrabCut-based methods which tried to increase the certainty

Table 1. Training configurations.

Backbone	U-Net	HNN [9]	ARU-Net [11]	Swin Transformer [22]
Pretrain	None	None	None	ADE20K [23]
Batch size	6	16	10	6
Prepare epochs	250	250	300	150
Total epochs	600	600	1000	600
Optimizer		AdaMax [24]		AdamW [25]
Learning rate	1×10^{-3}	1×10^{-3}	1×10^{-3}	6×10^{-5}
Scheduler	None	None	None	None
GPU	NVIDIA RTX 2080TI×1			
Augmentation	Flipping, cropping, padding, rotating, and scaling			

of the supervising mask via engineering with empirical rules. To this end, our work presents a new perspective of utilizing the RECIST annotation for weak supervision. It is worth noting that our label-space perturbation is also inspired by the data-space perturbations (i.e., data augmentation) commonly used in self-supervised contrastive learning [15].

3. EXPERIMENTS AND RESULTS

Dataset. We evaluate the proposed framework on the publicly available 2019 Kidney Tumor Segmentation (KiTS19) challenge data [20], which provides kidney tumor masks of 210 abdominal CT scans of unique patients in arterial phase. We randomly split the 210 volumes into training and test sets in the ratio of 80:20. For minimal preprocessing, the slices are rescaled to the range [0, 1] with a soft-tissue CT window range of [0, 400] Hounsfield unit (HU), and resized to 512×512 pixels. For weak supervision, we follow [9, 10] to convert the annotation masks to RECIST diameters by measuring the major and minor axes on 2D slices. The Dice score, Jaccard index, and 95th percentile of the Hausdorff distance (HD95) are used as evaluation metrics as in [21].

Implementation. All experiments are conducted with PyTorch 1.7.1 [26]. For backbone, we mainly consider the U-Net [27] and Swin Transformer [22] (base). The former is arguably the mostly widely used architecture for medical image segmentation, and the latter is a recent SOTA model for general image segmentation. We also implement our framework with two other backbones advocated in related works: the holistically nested networks (HNNs) [9] and ARU-Net [11], to demonstrate its model-agnostic applicability. As our purpose is to validate the effectiveness of the proposed dual-label consistency loss, no advanced post-processing trick is implemented; instead, we simply apply a default threshold of 0.5 to identify lesion pixels. The soft Dice loss [28] is used for $\ell_{\mathcal{P}}$ and $\ell_{\mathcal{A}}$. We first train the model with only the pseudo-label supervision by \mathcal{L}_{sup} until convergence (the “preparation”), and then add the consistency loss \mathcal{L}_{con} for remaining epochs. This is because \hat{Q} and \hat{C} should be reasonable under- and over-segmentations, respectively, for \mathcal{L}_{con} to contrast them validly. More details about the training configurations are charted in Table 1. Our codes will be released.

Dual Mask Validation. First, we verify the fundamental assumption that the constructed dual masks Q and C can serve as under- and over-segmenting supervision, respectively, by

Table 2. Lesion segmentation performance and comparison to SOTA approaches (mean±margin of error at 95% confidence level). The strong baseline [17] and upper bound are fully supervised models trained with GrabCut-generated and ground truth masks, respectively. A fixed $\lambda = 0.4$ is used for our framework, given its robustness to different λ values (cf. Fig. 3). *: significance at 0.05 level for pairwise comparison to our framework with the same backbone.

Method	Backbone	Dice \uparrow	Jaccard \uparrow	HD95 (pixel) \downarrow
Baseline	U-Net	0.854±0.025	0.778±0.028	4.083±0.204
WSSS [9]	HNN	0.838±0.030	0.766±0.032	4.118±0.233*
DRL [10]	U-Net+Deep Q-net	0.858±0.023	0.779±0.025	3.963±0.181
RLS [13]	U-Net	0.838±0.021*	0.744±0.023*	3.989±0.185
MULAN [18]	Mask R-CNN	0.878±0.022	0.808±0.025	3.718±0.199
SEENet [11]	Mask R-CNN+ARU-Net	0.890±0.020	0.825±0.023	3.516±0.191
Ours	HNN	0.845±0.028	0.770±0.029	3.889±0.180
	U-Net	0.862±0.026	0.792±0.028	3.955±0.206
Upper bound	ARU-Net	0.894±0.019	0.827±0.022	3.482±0.161
	U-Net	0.866±0.029	0.805±0.030	3.758±0.225*
Method	Backbone	Dice \uparrow	Jaccard \uparrow	HD95 (pixel) \downarrow
Baseline	Swin-T	0.891±0.018*	0.822±0.021*	3.449±0.144
RLS [13]	Swin-T	0.888±0.012*	0.808±0.015*	3.354±0.116
Ours	Swin-T	0.907±0.016	0.846±0.020	3.316±0.158
Upper bound	Swin-T	0.912±0.018	0.856±0.021	3.312±0.117

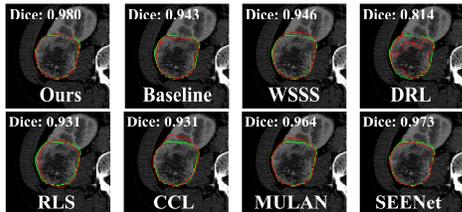


Fig. 2. Segmentation results of different methods (green: ground truth; red: network prediction; zoomed for better visibility). Baseline, RLS, and our framework use the Swin Transformer backbone, whereas the others use the backbone networks suggested in the original papers.

evaluating their recall and precision against the ground truth M . For ideal under segmentation, the recall should be low while the precision should be high, and vice versa for over segmentation. The recall and precision are 0.715 and 0.990 for Q , and 0.982 and 0.802 for C , verifying the assumption. Besides, compared to fitting an ellipse to the RECIST diameters [13] (recall 0.954 and precision 0.854), our C is a better over segmentation as indicated by its apparently higher recall and lower precision, thus more suitable for our framework.

Lesion Segmentation. We evaluate the segmentation performance of our proposed framework and compare to several SOTA approaches. The evaluation is done with whole-slice input (instead of cropped LOIs) for more practical use scenarios, albeit more challenging. For compared methods, we mainly use the backbones suggested in the original papers and follow the optimal training schedules described thereby; otherwise the U-Net and Swin Transformer are adopted and optimized for model-agnostic approaches. The results are shown in Table 2. As expected, the choice of the backbone network has a major impact on the performance, with the Swin Transformer performing the best. Despite that, when

Table 3. Ablation study (with U-Net) results in Dice scores (mean± margin of error at 95% confident level).

ℓ_P	ℓ_A	\hat{Q}	\hat{C}	$(\hat{Q} + \hat{C})/2$
×	×	0.683±0.023	0.809±0.024	0.826±0.023
✓	×	0.781±0.024	0.829±0.022	0.848±0.022
×	✓	0.844±0.022	0.852±0.021	0.862±0.026

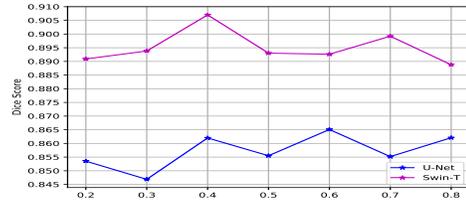


Fig. 3. Effects of varying λ (x axis) on performance (y axis).

the same backbones are used, our framework not only improves upon the strong baseline in all metrics, but also generally outperforms all competing methods. In fact, the Dice scores of our framework are fairly close to those of the upper bound (trained with pixel-wise full supervision by ground truth masks) using the same backbones. These results demonstrate the effectiveness and robustness of our proposed framework, although it is simpler in terms of both mask initialization (compared to the GrabCut) and network design (compared to the DRL [10] and SEENet [11]). Fig. 2 shows example segmentation results by different methods.

Ablation Study. We conduct ablation studies with and without co-training by \mathcal{L}_{con} , and with \mathcal{L}_{con} effective on the whole slice (ℓ_P) and in the constrained region (ℓ_A). As shown in Table 3, with ℓ_P added, substantial improvements are observed in both \hat{Q} and \hat{C} , as well as the ensemble; when additionally constraining the effective region of the consistency loss with ℓ_A , further improvements are achieved. Eventually, the ensemble of \hat{Q} and \hat{C} achieves the best performance with ℓ_A . These results demonstrate the efficacy of the proposed label-space co-training framework, especially with the region-constrained consistency loss. In addition, our framework introduces only one hyperparameter, i.e., λ in Eqn. (5) that controls the relative importance of the co-training loss. As shown in Fig. 3, the small variations (less than 0.025 Dice score) corresponding to both backbones demonstrate that our framework is not sensitive to the exact value of λ in a reasonable range ([0.4, 0.7]).

4. CONCLUSION

This work presented a novel co-training framework for lesion segmentation in CT slices, with weak supervision by the RECIST annotations, which effectively co-trained two sub-nets with a novel label-space perturbation induced consistency loss. Extensive experiments validated the efficacy of co-training with the proposed consistency loss, our framework’s superiority to existing works, and its model-agnostic property.

Compliance with Ethical Standards. This research study

was conducted retrospectively using human subject data made available in open access by the KiTS19 challenge [20]. Ethical approval was not required as confirmed by the license attached with the open access data.

Conflicts of Interest. The authors have no relevant financial or non-financial interests to disclose.

Acknowledgments. This work was supported by the Ministry of Science and Technology of the People’s Republic of China (STI2030-Major Projects2021ZD0201900).

5. REFERENCES

- [1] H. Beaumont et al., “Radiology workflow for RECIST assessment in clinical trials: Can we reconcile time-efficiency and quality?,” *Eur. J. Radiol.*, vol. 118, pp. 257–263, 2019.
- [2] E.A. Eisenhauer et al., “New response evaluation criteria in solid tumours: revised RECIST guideline (version 1.1),” *Eur. J. Cancer*, vol. 45, no. 2, pp. 228–247, 2009.
- [3] D. Wei et al., “Medical image segmentation and its application in cardiac MRI,” *Biomedical Image Understanding, Methods and Applications*, pp. 47–89, 2015.
- [4] M.H. Hesamian, W. Jia, X. He, and P. Kennedy, “Deep learning techniques for medical image segmentation: Achievements and challenges,” *J. Digit. Imaging*, vol. 32, no. 4, pp. 582–596, 2019.
- [5] S. Wang et al., “Conquering data variations in resolution: A slice-aware multi-branch decoder network,” *IEEE Trans. Med. Imaging*, vol. 39, no. 12, pp. 4174–4185, 2020.
- [6] J. Dai, K. He, and J. Sun, “BoxSup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation,” in *ICCV*, 2015, pp. 1635–1643.
- [7] D. Lin, J. Dai, J. Jia, K. He, and J. Sun, “ScribbleSup: Scribble-supervised convolutional networks for semantic segmentation,” in *CVPR*, 2016, pp. 3159–3167.
- [8] V. Agarwal, Y. Tang, J. Xiao, and R.M. Summers, “Weakly-supervised lesion segmentation on CT scans using co-segmentation,” in *Medical Imaging: Computer-Aided Diagnosis*. International Society for Optics and Photonics, 2020, vol. 11314.
- [9] J. Cai et al., “Accurate weakly-supervised deep lesion segmentation using large-scale clinical annotations: Slice-propagated 3D mask generation from 2D RECIST,” in *MICCAI*. Springer, 2018, pp. 396–404.
- [10] Z. Li and Y. Xia, “Deep reinforcement learning for weakly-supervised lymph node segmentation in CT images,” *IEEE J. Biomed. Health Inform.*, vol. 25, no. 3, pp. 774–783, 2020.
- [11] Y. Tang, K. Yan, J. Xiao, and R.M. Summers, “One click lesion RECIST measurement and segmentation on CT scans,” in *MICCAI*. Springer, 2020, pp. 573–583.
- [12] Y. Tang et al., “Lesion segmentation and RECIST diameter prediction via click-driven attention and dual-path connection,” in *MICCAI*. Springer, 2021, pp. 341–351.
- [13] Y. Tang et al., “Weakly-supervised universal lesion segmentation with regional level set loss,” in *MICCAI*. Springer, 2021, pp. 515–525.
- [14] C. Rother et al., “‘GrabCut’ interactive foreground extraction using iterated graph cuts,” *ACM Trans. Graph.*, vol. 23, no. 3, pp. 309–314, 2004.
- [15] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *CVPR*, 2020, pp. 9729–9738.
- [16] A. Blum and T. Mitchell, “Combining labeled and unlabeled data with co-training,” in *COLT*, 1998, pp. 92–100.
- [17] M. Zlocha, Q. Dou, and B. Glocker, “Improving RetinaNet for CT lesion detection with dense masks from weak RECIST labels,” in *MICCAI*. Springer, 2019, pp. 402–410.
- [18] K. Yan et al., “MULAN: Multitask universal lesion analysis network for joint lesion detection, tagging, and segmentation,” in *MICCAI*. Springer, 2019, pp. 194–202.
- [19] Y. Wang et al., “ACN: Adversarial co-training network for brain tumor segmentation with missing modalities,” in *MICCAI*. Springer, 2021, pp. 410–420.
- [20] N. Heller et al., “The KiTS19 challenge data: 300 kidney tumor cases with clinical context, CT semantic segmentations, and surgical outcomes,” *arXiv preprint arXiv:1904.00445*, 2019.
- [21] X. Luo, J. Chen, T. Song, and G. Wang, “Semi-supervised medical image segmentation through dual-task consistency,” *AAAI*, vol. 35, no. 10, pp. 8801–8809, 2021.
- [22] Z. Liu et al., “Swin transformer: Hierarchical vision transformer using shifted windows,” in *CVPR*, 2021, pp. 10012–10022.
- [23] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, “Scene parsing through ADE20K dataset,” in *CVPR*, 2017, pp. 633–641.
- [24] D.P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [25] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *ICLR*. 2019, OpenReview.net.
- [26] A. Paszke et al., “PyTorch: An imperative style, high-performance deep learning library,” *NeurIPS*, vol. 32, pp. 8026–8037, 2019.
- [27] O. Ronneberger, P. Fischer, and T. Brox, “U-Net: Convolutional networks for biomedical image segmentation,” in *MICCAI*. Springer, 2015, pp. 234–241.
- [28] F. Milletari, N. Navab, and S.-A. Ahmadi, “V-Net: Fully convolutional neural networks for volumetric medical image segmentation,” in *3DV*. IEEE, 2016, pp. 565–571.