

Unleashing the Strengths of Unlabeled Data in Pan-cancer Abdominal Organ Quantification: the FLARE22 Challenge

Jun Ma, Yao Zhang, Song Gu, Cheng Ge, Shihao Ma, Adamo Young, Cheng Zhu, Kangkang Meng, Xin Yang, Ziyang Huang, Fan Zhang, Wentao Liu, YuanKe Pan, Shoujin Huang, Jiacheng Wang, Mingze Sun, Weixin Xu, Dengqiang Jia, Jae Won Choi, Natália Alves, Bram de Wilde, Gregor Koehler, Yajun Wu, Manuel Wiesenfarth, Qiongjie Zhu, Guoqiang Dong, Jian He, the FLARE Challenge Consortium, and Bo Wang

Abstract

Quantitative organ assessment is an essential step in automated abdominal disease diagnosis and treatment planning. Artificial intelligence (AI) has shown great potential to automatize this process. However, most existing AI algorithms rely on many expert annotations and lack a comprehensive evaluation of accuracy and efficiency in real-world multinational settings. To overcome these limitations, we organized the FLARE 2022 Challenge, the largest abdominal organ analysis challenge to date, to benchmark fast, low-resource, accurate, annotation-efficient, and generalized AI algorithms. We constructed an intercontinental and multinational dataset from more than 50 medical groups, including Computed Tomography (CT) scans with different races, diseases, phases, and manufacturers. We independently validated that a set of AI algorithms achieved a median Dice Similarity Coefficient (DSC) of 90.0% by using 50 labeled scans and 2000 unlabeled scans, which can significantly reduce annotation requirements. The best-performing algorithms successfully generalized to holdout external validation sets, achieving a median DSC of 89.5%, 90.9%, and 88.3% on North American, European, and Asian cohorts, respectively. They also enabled automatic extraction of key organ biology features, which was labor-intensive with traditional manual measurements. This opens the potential to use unlabeled data to boost performance and alleviate annotation shortages for modern AI models.



INTRODUCTION

Abdominal organs are high cancer incidence areas, such as liver cancer, kidney cancer, pancreas cancer, and gastric cancer [1]. Computed Tomography (CT) scanning has been a major imaging technology for the diagnosis and treatment of abdominal cancer because it can yield important prognostic information with fast imaging speed for cancer patients, which has been recommended by many clinical treatment guidelines. In order to quantify abdominal organs, radiologists and clinicians need to manually delineate organ boundaries in each slice of the 3D CT scans [2], [3]. However, manual segmentation is time-consuming and inherently subjective with inter- and intra-expert variability. Inaccurate segmentation and quantification can lead to over-dosing or under-dosing treatment, increasing the risk of injuries. Therefore, accurate and automatic segmentation and quantification tools are highly demanded for abdominal cancer treatment in clinical practice.

AI has shown great promise for automatic cancer quantification and diagnosis in medical images [4]–[7]. In the field of abdominal cancer CT analysis, a large number of automatic segmentation algorithms have been proposed during the past decade, specifically for liver cancer [8], kidney cancer [9], and pancreas cancer [10]. However, these algorithms face limitations due to their reliance on costly annotations and lack of comprehensive evaluation on multinational datasets. AI competitions have been an effective way to gather efforts from the whole research community to solve hard tasks and promote methodology developments [11]–[13]. Several abdominal cancer-related AI challenges have been organized [14]–[16]. These challenges greatly facilitated innovations in liver and kidney cancer image analysis but the algorithms were not evaluated on external datasets and the evaluation metrics only focused on segmentation accuracy without considering algorithm efficiency. Therefore, there is an unmet need for annotation-efficient and universal abdominal organ segmentation algorithms that are robust to different diseases, races, and imaging platforms.

To address the above limitations, we organized a global AI challenge, the Fast and Low-resource semi-supervised Abdominal oRgan sEgmentation (FLARE), with diverse abdominal CT images and novel challenge design to prompt the development of annotation- and resource-efficient AI algorithms for abdominal organ quantification. Specifically, we constructed a large-scale multi-racial, multi-center, multi-disease, multi-phase, and multi-manufacturer abdominal CT

- A full list of affiliations appears at the end of the paper. Corresponding author: Bo Wang. E-mail: bowang@vectorinstitute.ai

dataset from 2900 patients, including 725,000 slices, 53 medical groups, seven CT scanner manufacturers, and four CT phases and covering more than six types of abdominal cancer. This is the largest and most diverse publicly available dataset to date. Furthermore, we designed a semi-supervised competition task, where participants were provided with a limited set of annotated images and a large number of unlabeled images. The evaluation metrics focus on both segmentation accuracy and efficiency. In order to validate generalization ability, top-performing algorithms were externally evaluated on independent Asian, European, and North American cohorts. The validation datasets were blind to participants, minimizing the risk of data leakage and providing a fair platform for benchmarking. This is the first time that AI algorithms for abdominal organ quantification have been challenged to learn with limited annotations and generalize on unseen medical centers without user interaction or additional fine-tuning. The dataset, top algorithms, and benchmark platform have been made publicly available for reproduction and long-term algorithm comparison.

RESULTS

Challenge design

The FLARE 2022 challenge followed two guidelines: Enhancing the QUALity and Transparency Of health Research (EQUATOR, <https://www.equator-network.org/>) and Biomedical Image Analysis ChallengeS (BIAS) [17], which has been pre-registered and passed peer review (two technical reviews and one clinical review) [18]. The challenge aims to benchmark algorithms that can automatically segment 13 abdominal organs from CT scans, including the liver, right kidney, spleen, pancreas, aorta, inferior vena cava, right adrenal gland, left adrenal gland, gallbladder, esophagus, stomach, duodenum, and left kidney (Fig. 1a). These organs were selected based on clinical needs and featured representative difficulties that are usually encountered in medical image analysis tasks.

Motivated by the real-world setting, we designed a semi-supervised segmentation task (Fig. 1b) rather than the typical fully supervised learning task because an abundance of unlabeled medical images exists in most medical centers, and collecting a small annotated dataset is feasible. Specifically, participants were provided with 2000 unlabeled cases and 50 labeled cases to develop automatic segmentation algorithms for abdominal organ segmentation. The challenge contained two phases. During the development phase, participants can access the whole training set and the 50 images in the tuning set. The tuning set labels were deployed on the online evaluation platform that was not publicly available, but participants can submit their results to the platform and get feedback on the algorithm performance.

During the validation phase, each team had one chance to submit one algorithm as the final solution. Different from most of the existing challenges [11], [13] that only considered accuracy-related metrics during ranking, our evaluation metrics focused on both segmentation accuracy and efficiency. In particular, we evaluated the degree of region overlap (Dice Similarity Coefficient, DSC), the degree of boundary matching (Normalized Surface Distance, NSD), running time, GPU memory consumption, and CPU utilization (Methods). This is the first time that medical image segmentation algorithms have been quantitatively evaluated with both accuracy-related metrics and comprehensive resource (i.e., GPU and CPU) consumption metrics. The challenge lasted 122 days and attracted 1,616 submissions from 112 participants on the tuning set evaluation platform, resulting in a total of 58,100 predictions. After the main challenge event at MICCAI, we collected three new cohorts and conducted the post-challenge analysis, aiming to evaluate the generalization ability of the top-performing algorithms.

Dataset characteristic

To create a comprehensive representation of real-world scenarios, we curated a robust and diverse dataset sourced from 53 distinct medical groups. This dataset encompassed a wide range of attributes, including varying racial backgrounds, multiple diseases, different imaging phases, and diverse CT machine manufacturers (Fig. 1c, Table 1, Supplementary Table 1-3). In total, 2900 images were retrospectively de-identified for the algorithm development and independent validation, covering North American and European patients with typical abdominal diseases: liver cancer, kidney cancer, spleen disease, and pancreas cancer. The internal validation set consists of 200 cases where 100 cases with the same disease type as the development set and 100 cases with new abdominal diseases, such as stomach cancer, sarcomas, colon cancer, ovarian cancer, and bladder diseases. The external validation sets are from three individual cohorts with North American, European, and Asian patients, respectively. Each external validation set contains 200 labeled cases from multiple medical centers. Annotation quality is highly important for model performance [19] and we employed a hierarchical human-in-the-loop pipeline to maintain consistent annotations (Methods).

Overview of evaluated algorithms

We received 48 successful algorithm docker containers during the validation phase. Each algorithm docker container was independently evaluated on the same workstation to obtain five quantitative metrics on the internal validation set and the final rank was generated with the ‘rank-then-aggregate’ strategy (Methods). The evaluation metrics and the ranking scheme has been made available to participants at the beginning of the challenge for the sake of fairness and transparency [20]. The top 10 teams were invited to join the FLARE consortium for post-challenge analysis and ablation study using unlabeled data. Four teams ranked top 10 in terms of segmentation accuracy but were not ranked top 10 overall because of low efficiency. However, in some clinical scenarios where algorithm efficiency is not the major consideration, such algorithms

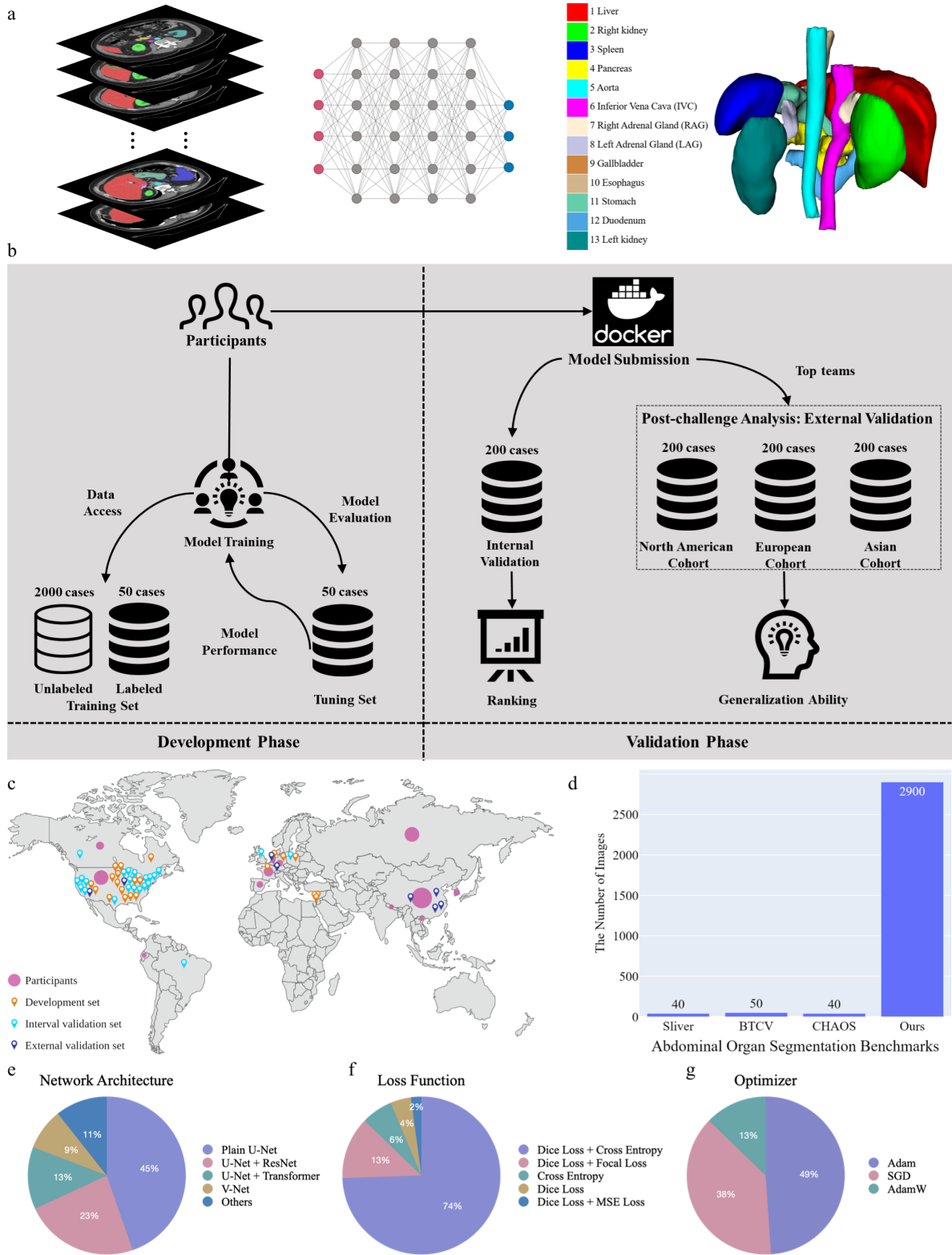


Fig. 1. Overview of the challenge design. **a**, The challenge aims to benchmark automatic algorithms that can simultaneously segment 13 abdominal organs. Organs have different sizes, morphologies, and appearances, featuring representative difficulties in medical image analysis tasks. **b**, The challenge contains two phases. During the development phase, participants develop automatic segmentation algorithms based on 2000 unlabeled cases and 50 labeled cases. The algorithms can be evaluated on the tuning set and the online evaluation platform will return the quantitative performance to participants. During the validation phase, each participant team can submit one algorithm via the docker container as the final solution, which is independently evaluated on the internal validation set to obtain ranking results. The top teams are selected for post-challenge analysis, which are further evaluated on three independent intercontinental cohorts to validate their generalization ability. **c**, The data sources are multinational and the challenge has attracted more than 100 worldwide participants (the circle size is proportional to the number of participants in each country). **d**, The FLARE challenge dataset is significantly larger than the previous abdominal organ segmentation challenge datasets. Distribution of key algorithm designs: **e** network architecture, **f**, loss function, and **g**, optimizer.

TABLE 1

Dataset characteristics of the development set, internal validation set, and three external validation sets. The development set includes a training set with 50 labeled images and 2000 unlabeled images and a tuning set with 50 images. All the images in the development set are available to participants but the tuning set annotations are not released. All validation sets are fully independent and hidden from participants. For each organ, we report the statistics of volume or diameter. Values are displayed as Median values (First quartile, Third quartile).

Name	Development set	Internal validation set	North American external validation set	European external validation set	Asian external validation set
No. of cases	2100	200	200	200	200
No. of sources	22	23	2	2	4
Region	North American European	North American European	North American	European	Asian
Disease	Liver cancer Kidney cancer Spleen disease Pancreas cancer	Liver cancer Kidney cancer Spleen disease Pancreas cancer Stomach cancer Sarcomas Ovarian cancer Bladder disease	Liver cancer Pancreas cancer	Various abdominal disease	Various abdominal disease
Phase	Plain phase Artery phase Portal phase Delay phase	Plain phase Artery phase Portal phase Delay phase	Plain phase Artery phase Portal phase Delay phase	Plain phase Artery phase Portal phase Delay phase	Plain phase Artery phase Portal phase Delay phase
Manufacturer	Siemens, GE, Philips, Toshiba, Barco, Vital	Siemens, GE, Philips, Toshiba, Imatron, Vital PHMS	GE, Philips, Toshiba	Siemens, GE, Philips, Toshiba,	Siemens, GE, Philips, Toshiba,
Liver volume	1601 (1383, 1921)	1542 (1323, 1655)	1611 (1264, 1970)	1458 (1223, 1662)	1188 (1013, 1367)
Right kidney volume	197.6 (162.4, 254.5)	179.0 (152.5, 207.5)	170.6 (149.4, 212.7)	158.7 (125.2, 196.1)	141.8 (122.6, 161.9)
Left kidney volume	196.6 (166.8, 251.1)	184.4 (146.8, 213.1)	175.8 (153.6, 201.0)	153.7 (126.3, 193.1)	146.7 (126.5, 174.9)
Spleen volume	230.7 (158.8, 320.3)	230.3 (174.9, 308.4)	292.2 (164.4, 497.3)	182.3 (133, 271.4)	182.9 (130.0, 252.4)
Pancreas volume	89.16 (67.96, 119.3)	77.38 (64.34, 92.93)	76.14 (59.17, 90.09)	73.19 (58.55, 92.61)	74.26 (60.43, 91.70)
Aorta diameter	21.83 (20.40, 24.45)	23.24 (21.53, 25.44)	23.66 (21.30, 25.85)	24.32 (21.41, 27.08)	21.62 (19.42, 24.04)
Inferior vena cava	26.12 (24.43, 27.90)	26.01 (23.56, 27.84)	25.19 (23.63, 27.41)	25.83 (23.67, 27.71)	24.79 (23.03, 26.39)
Right adrenal gland	4.222 (3.551, 5.540)	4.802 (3.794, 5.573)	4.340 (3.350, 5.339)	3.829 (2.747, 4.994)	3.297 (2.513, 4.342)
Left adrenal gland	5.265 (4.348, 6.404)	5.423 (4.342, 6.474)	4.963 (3.938, 6.137)	4.195 (3.012, 5.564)	3.977 (3.187, 5.048)
Gallbladder volume	26.80 (16.92, 44.33)	28.94 (17.66, 41.78)	23.53 (14.04, 44.12)	22.5 (12.08, 34.27)	18.17 (8.734, 27.76)
Esophagus diameter	19.11 (17.16, 21.66)	18.55 (16.70, 19.87)	20.49 (18.09, 22.72)	18.93 (16.54, 21.13)	17.85 (15.83, 20.57)
Stomach volume	308.7 (226.9, 419.6)	349.1 (256.3, 525.4)	371.2 (274.1, 510.5)	310.7 (208.2, 472.9)	342.3 (214.8, 502.2)
Duodenum diameter	76.52 (65.93, 87.74)	71.20 (61.71, 89.64)	74.32 (64.90, 88.99)	65.61 (52.55, 81.82)	58.23 (48.28, 72.93)

are also useful. Therefore, we also invited these teams to join the FLARE consortium and conducted the post-challenge analysis. We refer to the algorithms of the 14 teams as the top-performing algorithms.

We analyzed the characteristics of the employed algorithms by participants (Fig. 1e-g, Table 4-5). All the submitted algorithms were based on deep learning. 80% teams used U-Net as the main network architecture (Fig. 1e), where 45% teams used plain U-Net and others combined U-Net with popular networks, such as ResNet (23%) and Transformer [21]. All the algorithms used dice loss [22] and its combination with cross-entropy loss was the most popular choice (Fig. 1f) because compound loss function has been proven to be robust [23]. Adam [24] was the most frequently used optimizer followed by the stochastic gradient descent (SGD) (Fig. 1g).

Best-performing algorithm. Team blackbean (T1 [25]) proposed a pseudo-labelling framework with nnU-Net [26] as the base network architecture. Multiple big nnU-Net models were trained based on the labeled cases. Then, unlabeled cases were passed to the trained models to predict pseudo-labels. After that, the big nnU-Net models were updated based on the combination of labeled cases and unlabeled cases with pseudo-labels. The process was iterated for three rounds and low-quality pseudo-labels were filtered by the uncertainty score which was defined based on the changes in pseudo-labels during different iterations. To improve the inference efficiency, a small nnU-Net was trained by reducing the input size and network size. Furthermore, a modified sliding window strategy was designed to filter background regions based on anatomy prior, which can reduce the computational cost.

Best-segmentation-accuracy algorithm. Team aladdin5 (T23 [27]) developed a cascaded framework with nnU-Net [26] and pseudo-label learning. A binary U-Net segmentation network was first trained to localize the ROI in low-resolution images. Then, an augmented U-Net with more trainable parameters was used to segment the 13 target organs from the ROI in high-resolution images. In addition, an ensemble of multiple U-Nets was used to segment the unlabeled images to generate pseudo-labels. Finally, the cascaded framework was trained with 50 labeled cases and 2000 cases with pseudo-labels. However, this framework is time-consuming because of the large model size. The inference speed was nearly 10 times slower than the top three best-performing teams.

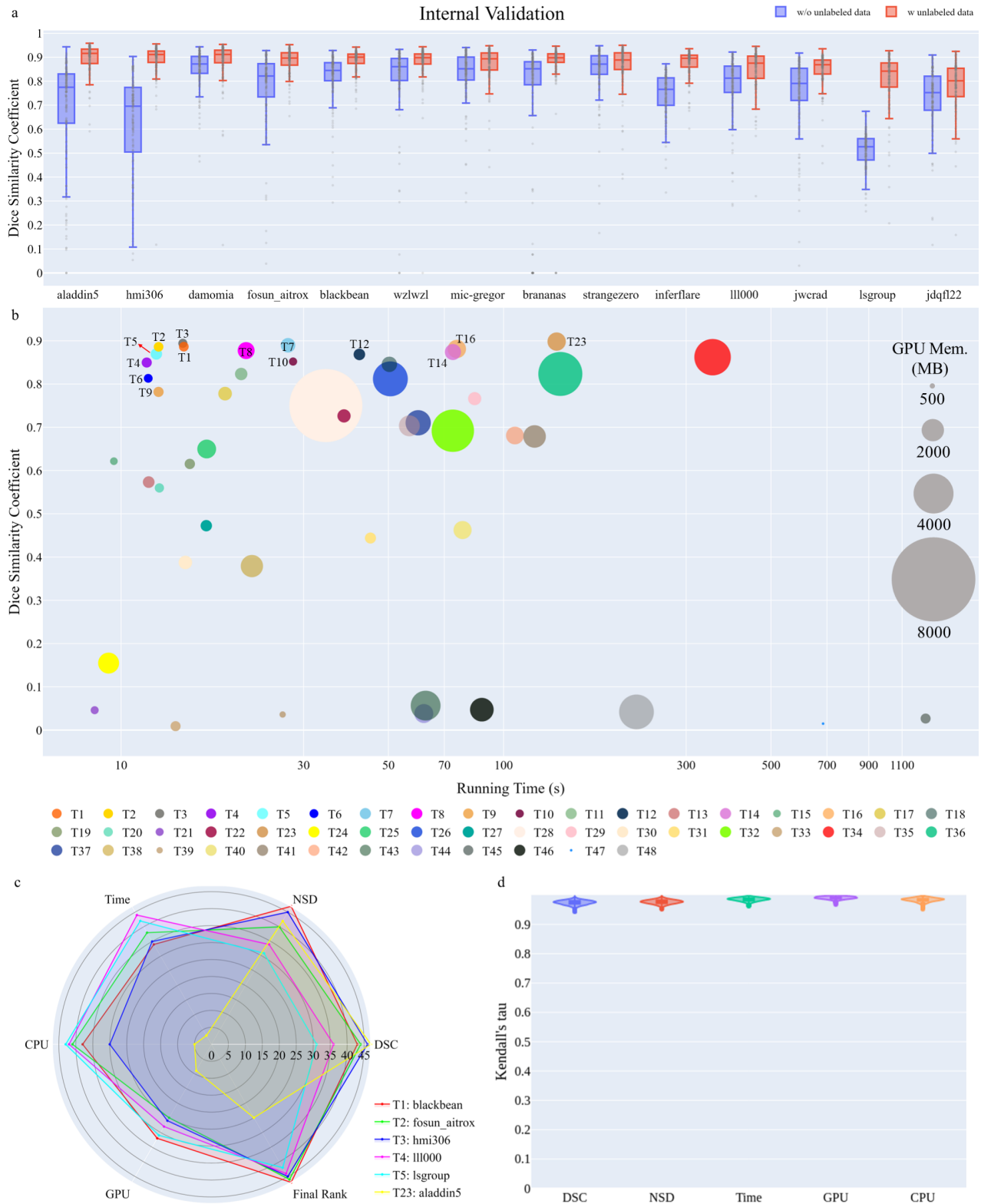


Fig. 2. **Performance analysis on the internal validation set.** **a**, The comparisons of using unlabeled data and without using unlabeled data for top-performing algorithms show that using unlabeled data can significantly improve the performance. The average improvement of the Dice Similarity Coefficient (DSC) score is 9.8%. **b**, The top three best-performing algorithms achieve a good trade-off between segmentation accuracy (y-axis) and efficiency (x-axis). The circle size is proportional to GPU memory consumption. The 14 top-performing algorithms are marked in the figure. **c**, The performance comparisons among different dimensions are presented for the top three best-performing algorithms and another three top-performing algorithms with the best DSC, running time, and CPU utilization metrics, respectively. The value denotes the number of teams surpassed by each algorithm in each dimension. **d**, The bootstrap distribution of rankings (N=1000) shows that the ranking scheme is stable with respect to sampling variability.

Performance analysis on the internal validation dataset

We independently evaluated all the submitted algorithms on the holdout internal validation set during the validation phase. Dice Similarity Coefficient (DSC) and Normalized Surface Distance (NSD) [28] were used for quantitative evaluation. As shown in Fig. 2a and Supplementary Table 6, the algorithms with the best DSC scores demonstrated remarkable agreement with the organ contours of the ground truth. The median average organ DSC score is 91.6% (interquartile range (IQR): 87.4-93.4%), significantly surpassing the performance of other teams (Wilcoxon signed-rank test, $p < 0.05$).

To evaluate the impact of unlabeled data, we conducted a comparison between the algorithms utilizing unlabeled data and their counterparts that did not employ the unlabeled data (Fig. 2a). Notably, all algorithms exhibited significant performance improvements (paired T-test $p < 0.05$) when leveraging unlabeled images. The team hmi306 achieved the highest performance gain of 26.8%. Furthermore, the incorporation of unlabeled cases contributed to a reduction in the number of cases with poor segmentation. These compelling results demonstrate that the use of unlabeled data has the potential to enhance model performance and improve generalization capabilities.

Algorithm efficiency is another important metric to consider when deploying models. We quantitatively evaluated the running time and resource usage of all algorithms (Fig. 2b, Methods, Supplementary Table 6, Fig. 2). For an input image with more than twenty million voxels ($\sim 512 \times 512 \times 100$), all the top ten algorithms (T1-T10) can finish the segmentation of 13 organs within 30 seconds. As depicted in the bubble plot (Fig. 2b), the majority of top-performing teams aimed to strike a balance between computing resource consumption and accuracy, evidenced by their concentration towards the upper left of the figure. Remarkably, the top three teams exhibited an excellent trade-off between segmentation accuracy and efficiency, achieving a median DSC of 88.6-89.4% in under 15 seconds, with average GPU memory usage below 2GB and CPU utilization below 30%. These findings highlight the notable efficiency achieved by the leading algorithms, making them well-suited for practical model deployments.

We further analyzed the characteristics exhibited by the top-performing algorithms with a radar plot that captured six dimensions for each algorithm (Fig. 2c). It revealed that algorithms solely prioritizing a single metric, such as DSC, failed to achieve top rankings. In contrast, algorithms striking a favorable trade-off between segmentation accuracy and efficiency secured prominent positions in the rankings. We also evaluated the ranking stability by performing bootstrap (1000 times) for all algorithms and computed Kendall's τ as a quantitative metric (Fig. 2d, Methods). For each metric, Kendall's τ scores were found to be very close to 1, indicating a high degree of agreement between the rankings. Additionally, the compact distributions of these scores further confirm the stability of the ranking results with respect to sampling variability. These findings provide robust evidence that the obtained rankings are highly consistent and reliable across different samples.

Finally, we conducted an in-depth analysis of the organ-wise segmentation performance exhibited by the 14 top-performing algorithms (Fig. 3). Notably, all the top algorithms achieved high DSC scores of over 90% for organs such as the liver, kidneys, and spleen, owing to their relatively larger volumes and simple shapes, which make them easier to segment accurately. For example, the best-segmentation-accuracy team aladdin5 achieved median DSC scores of 98.6% (IQR: 98.3-98.8%), 98.3% (IQR: 97.3-98.7%), 98.2% (IQR: 97.0-98.5%), and 98.6% (IQR: 98.0-99.0%) for the liver, right kidney, left kidney, and spleen, respectively.

In the case of tubular organs, the aorta and inferior vena cava (IVC) showed higher medical DSC scores than the esophagus. This discrepancy can be attributed to the lower contrast and thinner caliber of the esophagus, making it more challenging to segment. The segmentation of the remaining organs poses the greatest challenge. Specifically, the pancreas, stomach, and duodenum exhibit complex structures with boundaries that may overlap with surrounding organs. The adrenal glands and gallbladder, being relatively small within the abdomen, are also difficult to accurately recognize. These organs may be surgically removed in some cases, further complicating the segmentation task. As a result, the algorithms showcased varied performance when it came to these challenging organs. Notably, team aladdin5 achieved the best accuracy for these difficult organs, with median DSC scores ranging from 83.7% to 94.4%. Altogether, the top-performing algorithms demonstrated high segmentation accuracy for organs such as the liver, kidneys, and spleen, while facing challenges with more complex structures, smaller organs, and surgically altered organs, where diverse performance was observed.

Performance analysis on the external validation cohorts

One of the primary challenges limiting the widespread application of AI in clinical routines is the lack of generalization capability to new cohort data. To quantitatively assess the generalization performance, we conducted an independent evaluation of the top 14 algorithms on three external validation cohorts: the North American, European, and Asian cohorts (Fig. 4), which are curated from external medical centers (Supplementary Table 3).

The best-accuracy algorithm (aladdin5) generalized well on the external validation cohorts (Fig. 4a, Supplementary Fig. 3b, Table 6-8). In comparison to the internal validation set, the algorithm achieved comparable performance with median DSC scores of 89.3% (IQR: 84.4-93.0%), 90.9% (IQR: 84.3-94.2%), and 87.5 (IQR: 80.3-92.9%) on the North American, European, and Asian cohorts, respectively. Among the 14 top algorithms, multiple algorithms obtained competitive performance compared to the best DSC score, indicating their strong generalization ability on the external validation sets. Furthermore, organ-wise analysis (Supplementary Fig. 4-9) shows that these algorithms can accurately segment most organs (e.g., liver, kidneys, spleen, and aorta) with median DSC and NSD scores of more than 90%. However, for some challenging cases, such as organs with low contrast or severe pathological changes, there is still large room for further

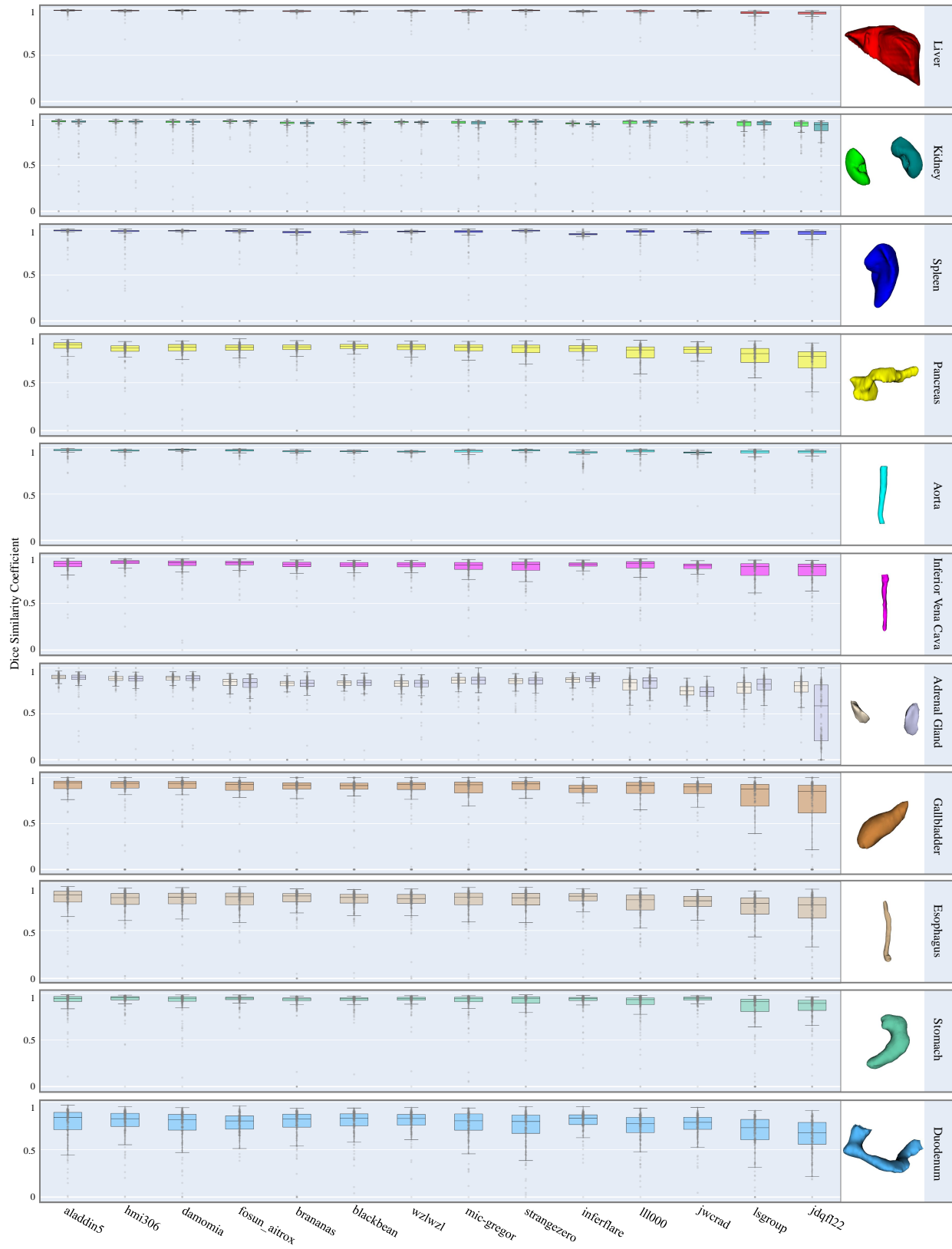


Fig. 3. Dot and box plots of the Dice Similarity Coefficient (DSC) values of top-performing algorithms for the 13 organs on the interval validation set. The box plots display descriptive statistics across all internal validation cases, with the median value represented by the black horizontal line within the box, the lower and upper quartiles delineating the borders of the box, and the vertical black lines indicating the 1.5 interquartile range. The algorithms are ranked on the x-axis based on their median DSC scores.

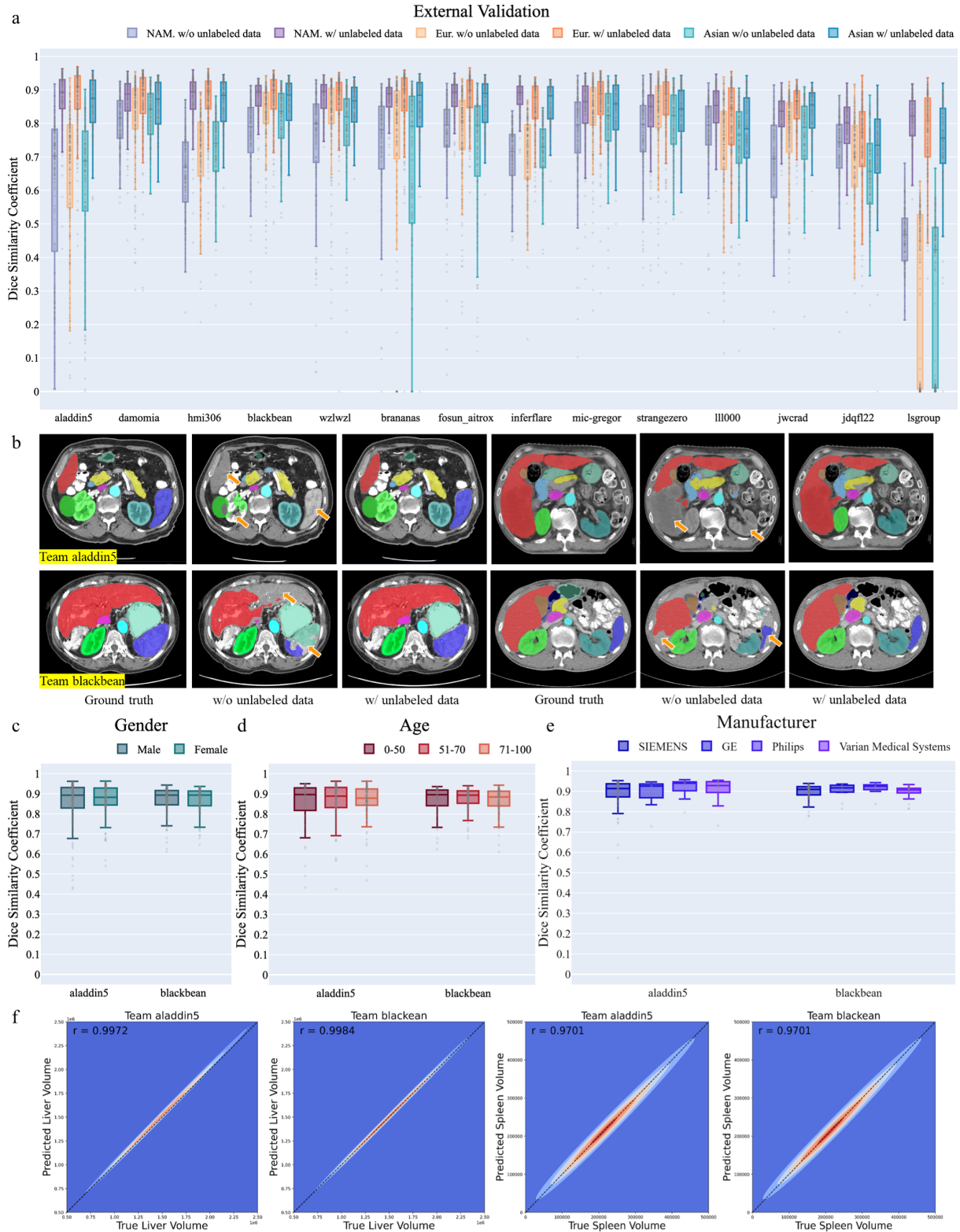


Fig. 4. **Performance on three external validation sets.** **a**, The segmentation performance (Dice similarity coefficient, DSC) on the North American (NAM.) cohort, European (Eur.) cohort, and Asian cohort. For each cohort, the DSC scores between using unlabeled data and without using unlabeled data are presented as well. **b**, Visualized segmentation examples of the two top algorithms show that using unlabeled data can significantly improve the segmentation quality. **c-e**, The segmentation performance of the best-accuracy algorithm (aladdin5) and the best-performing algorithm (blackbean) across demographics, including genders, ages, and manufacturers. **f**, Pearson's correlation contour plots of the organ volume demonstrate that the two top algorithms accurately quantify the liver and spleen volumes, which are important clinical biomarkers.

improvements (Fig. 10-12). Moreover, compared to the counterpart algorithms without unlabeled data, the top algorithms using unlabeled data still consistently achieved remarkable performance gains on the three external cohorts, indicating that using the unlabeled data during the development phase can significantly improve the model generalization ability.

In order to understand how the unlabeled data improves the segmentation quality, we visualized typical segmentation examples from the best-accuracy algorithm (T23: aladdin5) and the best-performing algorithm (T1: blackbean) (Fig. 4b, Supplementary Fig. 13-15). It can be found that the baseline model (without using unlabeled data) is not robust to different image contrasts, generating various over-segmentation and under-segmentation errors. In contrast, these segmentation errors were significantly reduced by learning with unlabeled data, demonstrating that unlabeled data can be used to improve the model generalization ability.

AI models have been proven that they could have biases on different genders, ages, and other demographic factors [29]. Therefore, it is necessary to evaluate the fairness of the top algorithms. We evaluated the fine-grained performance of the best-accuracy algorithm (T23: aladdin5) and the best-performing algorithm (T1: blackbean) across different genders, ages, and manufacturers on the external validation set. We separate gender groups to male and female, age groups to 0-50, 51-70, and over 70, manufacturer groups to SIEMENS, General Electric (GE), Philips, and Varian Medical System. The results imply that the two algorithms do not significantly vary regarding gender (Fig. 4c and Supplementary Fig. 3c), age (Fig. 4d and Supplementary Fig. 3d), and manufacturer (Fig. 4e and Supplementary Fig. 3e), highlighting their potential for broad applications in abdominal CT segmentation tasks. Moreover, the resulting algorithms enable automatic and high-throughput image-based phenotyping. We computed liver and spleen volume based on the segmentation results of algorithms aladdin5 and blackbean (Fig. 4f). There is a strong correlation (Pearson's $r=0.9701-0.9984$) between the predicted and ground-truth volume. This indicates that the two top algorithms accurately quantify organ volumes, which could replace manual linear measurements.

DISCUSSION

The FLARE 2022 challenge presented the largest challenge in abdomen image analysis in terms of both dataset size and the number of docker-packed algorithm submissions. The datasets covered real-world diversities of races, patients, abdominal cancer types, imaging manufacturers, and imaging protocols. The challenge attracted participants from all over the world with different research backgrounds, such as computer vision, biomedical engineering, biology, and radiology. Moreover, this was the first documented attempt to analyze the usage of unlabeled data for boosting segmentation performance and quantitative evaluation of algorithm efficiency and resource consumption. The resulting algorithms may reduce the time to extract abdominal organ biomarkers, empowering radiologists and clinicians to adopt and conduct more quantitative analysis for research and clinical practice.

There are three main findings based on the validation results. First, unlabeled data can be used to alleviate the annotation shortage problem and significantly improve algorithm generalization ability. Most of the participants employed similar network architectures: U-Net-like fully convolutional networks. The main factor to distinguish the best-performing teams was the successful exploitation of unlabeled data with pseudo-label learning. Second, the results demonstrate that AI algorithms can achieve a great trade-off between accuracy, efficiency, and resource consumption. Third, most organs have obtained very high agreements with reference standards but some organs still have a large room for further improvements, such as the pancreas, stomach, and duodenum because of their irregular shapes and dynamic appearances. This indicates that abdominal organ segmentation is still an unsolved problem (Supplementary Fig. 10-12).

We also identify some useful strategies on how to use unlabeled data to improve algorithm performance. All the best-performing teams used pseudo-label learning to incorporate the unlabeled data into the algorithm development. The main idea is to train models on a small well-labeled dataset and predict the unlabeled data to generate pseudo-labels. Then, one can update the model weights with all the data. In this pipeline, the key is to generate high-quality pseudo-labels. Participants have developed different strategies to improve the quality of pseudo-labels. For example, team aladdin5 generated pseudo-labels by model ensembles which can incorporate knowledge from different models, while team blackbean defined an uncertainty metric to select confident pseudo-labels. By contrast, team fosun_aitrox employed an easy-to-hard curriculum learning manner to prioritize reliable pseudo-labels. The resulting algorithms have been extensively validated on the internal and external validation sets and demonstrated their effectiveness to bring remarkable performance gains. In addition, data augmentation was also used by all top teams, such as rotation, scaling, and intensity shifting. This augmentation transforms can improve algorithms' robustness to different CT imaging protocols.

Resource consumption and efficiency are also important factors during algorithm deployments in clinical environments because most medical centers do not have powerful computing devices. Many top teams employed two-stage frameworks. In the first stage, a small network was used roughly locate the abdomen region of interest (ROI). In the second stage, a larger network was used for fine-grained segmentation. To further improve the inference speed, one can resize the whole CT scan to a fixed image size in the first stage, which can reduce the number of voxels by 20 times. Moreover, image priors can also be used to reduce the computational burden. For example, the human tissues are located in the middle of the CT scans and the background intensities are zero. Thus, the algorithm prediction process can directly ignore the background and only predict the middle part of the input CT image.

This study has two main limitations. First, although we have constructed the largest abdomen CT dataset, the included patients were predominantly North American, European, and Asian, which lack patients from Africa. Second, this study

mainly focuses on organ quantification while tumor quantification is also important in clinical practice. In the near future, we will increase the tumor annotations for this dataset and extend the challenge with joint organ and tumor segmentation.

This paper presents a successful proof of concept that unlabeled data are indeed useful in AI algorithm development when annotations are in shortage. The top algorithms obtained consistent and significant improvement by using the unlabeled data. Moreover, abdominal organ segmentation is still an unsolved problem, especially for organs with complex shapes and diverse appearances. We anticipate that in the near future, radiologists and clinicians can be assisted by AI algorithms to objectively and automatically assess organ status from abdominal CT scans at high throughput. In addition, we have made this large abdominal dataset and the code repositories of top algorithms publicly available to the community for research use <https://flare22.grand-challenge.org/>. This open-source commitment not only facilitates reproducibility but also encourages collaboration and fosters progress in abdominal disease research.

FLARE challenge consortium

Junjun He, Hua Yang, Huihua Yang, Bingding Huang, Mengye Lyu, Yongkang Ma, Heng Guo, Rongguo Zhang, and Klaus Maier-Hein

Jun Ma is with the Department of Laboratory Medicine and Pathobiology, University of Toronto; Peter Munk Cardiac Centre, University Health Network; Vector Institute, Toronto, Canada

Yao Zhang is with Shanghai AI Laboratory, 200232, Shanghai, China

Song Gu is with the Department of Image Reconstruction, Nanjing Anke Medical Technology Co., Ltd., 211113, Nanjing, China

Cheng Ge is with Ocean University of China, 266100, Qingdao, China

Shihao Ma and Adamo Young are with the Department of Computer Science, University of Toronto; Peter Munk Cardiac Centre, University Health Network; Vector Institute, Toronto, Canada

Cheng Zhu is with Tinavi Medical Technologies Co., Ltd., 100192, Beijing, China

Kangkang Meng is with University of Science and Technology Beijing, 100083, Beijing, China

Xin Yang is with the School of Biomedical Engineering, Health Science Center, Shenzhen University, 518055, Shenzhen, China

Ziyan Huang is with Shanghai AI Laboratory, 200232, Shanghai, China

Fan Zhang is with the Department of Radiological Algorithm, Fosun Aitrox Information Technology Co., Ltd., 200033, Shanghai, China

Wentao Liu is with the School of Artificial Intelligence, Beijing University of Posts and Telecommunications, 100876, Beijing, China

Yuanke Pan and Shoujin Huang are with Shenzhen Technology University, 518000, Shenzhen, China

Jiacheng Wang is with Xiamen University, 361005, Xiamen, China

Mingze Sun is with Alibaba, 100084, Beijing, China

Weixin Xu is with Infervision Medical Technology Co., Ltd., 100025, Beijing, China

Dengqiang Jia is with Hong Kong Centre for Cerebro-cardiovascular Health Engineering, 000000, Hong Kong, China

Jae Won Choi is with the Department of Radiology, Armed Forces Yangju Hospital, 11429, Yangju, Korea

Natália Alves and Bram de Wilde are with the Department of Radiology, Radboudumc, 6525XZ, Nijmegen, Netherlands

Gregor Koehler is with the Department of Medical Image Computing, German Cancer Research Center, 69120, Heidelberg, Germany

Yajun Wu is with ShenZhen Yorktal DMIT LLC, 518100, Shenzhen, China

Manuel Wiesenfarth is with the Division of Biostatistics, German Cancer Research Center (DKFZ), Heidelberg 69120, Germany

Qiongjie Zhu, Guoqiang Dong, and Jian He are with the Department of Nuclear Medicine, Nanjing Drum Tower Hospital, 210008, Nanjing, China

Junjun He is with Shanghai AI Laboratory, 200232, Shanghai, China

Hua Yang is with the Department of Radiological Algorithm, Fosun Aitrox Information Technology Co., Ltd., Shanghai, China

Huihua Yang is with the School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, 100876, China

Bingding Huang is with the College of Big Data and Internet, Shenzhen Technology University, Shenzhen, 518188, China

Mengye Lyu is with the College of Health Science and Environmental Engineering, Shenzhen Technology University, Shenzhen, China

Yongkang Ma is with the Manteia Technologies Co.,Ltd, Xiamen, China

Heng Guo is with the Alibaba DAMO Academy, Beijing, China

Rongguo Zhang is with Infervision Medical Technology Co., Ltd., 100025, Beijing, China

Klaus Maier-Hein is with the Department of Medical Image Computing, German Cancer Research Center, 69120, Heidelberg, Germany

Bo Wang is with the Peter Munk Cardiac Centre, University Health Network; Department of Laboratory Medicine and Pathobiology and Department of Computer Science, University of Toronto; Vector Institute, Toronto, Canada

METHODS

Study design

The study design of the FLARE 2022 challenge was preregistered [18] and passed peer review on the 25th International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI 2022). The dataset was retrospectively obtained from 53 different medical sources (Supplementary Table 1-3). During the development phase, participants could access the 50 labeled cases from one medical center and 2000 unlabeled cases from 21 medical centers. All the internal and external validation sets were hidden from participants. For fair comparisons, participants were not allowed to use additional data or pretrained models to develop their methods.

We launched the challenge on 15 March 2022 via the GrandChallenge platform. During the development phase, each team had three chances per day to make submissions to the online evaluation platform and get segmentation accuracy scores. Moreover, each team also had five chances to submit their algorithm docker containers to challenge organizers and obtain segmentation efficiency scores. During the validation phase, each team was required to submit the final algorithm by 15 July 2022. The algorithm should be packaged as a docker container and we independently evaluated the algorithm docker on the internal validation set that is blind to participants.

Data standardization

All the images were de-identified and can be used in this challenge based on legal licenses (e.g., CC-BY, CC-BY-NC-SA). We normalized all the images to standard NIfTI format (<https://nifti.nimh.nih.gov/>) because it only contains necessary image information (intensity, orientation, spacing, origin) and without any patient metadata (e.g., age, gender). Compared to the well-known DICOM format (<https://www.dicomstandard.org/>), NIfTI aggregates multiple 2D slices to one 3D volume, which is easier to handle for developing AI algorithms. NIfTI has been the most popular data format in 3D medical image analysis challenges [12], [30]. The voxel orientation was standardized as canonical 'RAS', which means that the first, second, and third voxel axes go from left to Right, posterior to Anterior, and inferior to Superior, respectively.

Annotation protocols

The annotation process involved a team of five junior radiologists with one to five years of experience and two senior radiologists with more than 10 years of experience. The organs vary in different sizes, morphologies, and appearances. For example, the average liver volume (1500cm^3) is around 300 times larger than the volume of the adrenal gland (5cm^3). The esophagus, aorta, and inferior vena cava emerge as lumen structures, while the spleen and pancreas appear as irregular masses of tissues. Moreover, some organs are close to each other with low contrasts, such as the liver and stomach. Algorithms need to address all the difficulties at the same time since they occur simultaneously. Annotation quality is crucial to algorithm performance and it is necessary to reduce the annotation variability [19] between annotators. Therefore, we provided detailed annotation instructions to the annotators and employed a hierarchical human-in-the-loop pipeline to enhance throughput and maintain consistent annotations.

The hierarchical annotation pipeline consisted of three stages to ensure accuracy and consistency throughout the annotation process. In the first stage, an annotation consensus was formulated by a senior radiologist, an oncologist, and a surgeon based on radiation therapy oncology group consensus (RTOG) panel guideline [31] and Netter's anatomical atlas [32]. Specifically, the liver contour should include all hepatic parenchyma and all liver lesions. The hepatic vessels inside the liver also need to be covered. If the vessels are located outside the liver (i.e., the entrance of portal hepatitis) based on the coronal view, they should be excluded from the liver contour. The kidney contour should include the renal parenchyma while excluding adjacent structures such as blood vessels and surrounding fat. The spleen contour should include all splenic parenchyma and any splenic lesions. It should exclude adjacent structures such as the splenic vessels (arteries and veins), particularly those located outside the spleen. The pancreatic contour should encompass all pancreatic parenchyma including the head, body, and tail, as well as any pancreatic lesions. Exocrine, endocrine components, and the pancreatic duct need to be included, but the surrounding vessels and fat should be excluded. The aortic contour should include the entire lumen of the aorta, from the aortic root to the bifurcation. The aortic wall (including the aortic calcification) should also be included. The inferior vena cava contour should include the entire lumen and cover the walls. The adrenal gland contour should include the entire adrenal gland, both cortex and medulla, and any adrenal lesions. The gallbladder contour should encompass the entire gallbladder wall, including the body, fundus, and neck, as well as any gallstones or polyps. The cystic duct and the surrounding liver parenchyma should be excluded. The esophagus contour should include the entire esophageal wall, while adjacent structures such as the trachea, aorta, and surrounding fat and muscle should be excluded. The stomach contour should encompass the entire stomach wall including the fundus, body, antrum, and pylorus, as well as any gastric lesions. The duodenum contour should include the entire duodenal wall from the duodenal bulb to the ligament of Treitz, along with any duodenal lesions. It should exclude surrounding structures such as the head of the pancreas, common bile duct, and surrounding vasculature. Before the annotation process, all annotators were required to learn and follow the annotation consensus.

The second stage of our annotation pipeline involved a human-in-the-loop approach aimed at enhancing annotation throughput. To facilitate this process, we utilized five 3D U-Net models [26] trained via 5-fold cross-validation on existing abdomen CT datasets [33], [34]. Leveraging the predictions generated by these models, junior annotators performed manual

refinements on 100 randomly selected predictions, which were then checked and revised by the senior radiologists. This process was iterated seven times until all the images were labeled by one of the junior annotators and further refined by the senior radiologists. In the external validation set (600 CT scans), the organ annotations of 512 CT scans were generated by us, and the remaining 84 CT scans were collected from recent public abdomen datasets [35]–[37]. Their annotations were used and manually refined to follow our annotation consensus. It is important to note that these publicly available annotations were accessible only after the challenge, thereby preventing participants from utilizing them during the challenge. In the third and final stage, we aimed to identify potential annotation errors. We trained a new set of U-Net models using five-fold cross-validation, with special attention given to images exhibiting low Dice Similarity Coefficient (DSC) scores (≤ 0.75), which were then double-checked by the senior radiologists.

Evaluation platform

All the submitted algorithms were evaluated on the same workstation. Specifically, the workstation is a Ubuntu 20.04 desktop with one central processing unit (CPU, Intel Xeon(R) W-2133 CPU, 3.60GHz \times 12 cores), one graph processing unit (GPU, NVIDIA QUADRO RTX5000, 16G), 32G of memory, and 500G of hard disk drive storage. At the beginning of the challenge, we also released the versions of GPU Driver (510.60.02), CUDA (11.6), and Docker (20.10.13) to make sure that participants’ algorithms are compatible with our evaluation platform.

Evaluation metrics

We used five complementary metrics to quantitatively evaluate the segmentation accuracy, efficiency, and resource consumption of the algorithms (Supplementary). Specifically, Dice Similarity Coefficient (DSC), the most popular segmentation metric [20], was used to evaluate the spatial overlap between the segmentation mask and ground truth, which is defined by

$$DSC(G, S) = \frac{2|G \cap S|}{|G| + |S|},$$

where $|\cdot|$ counts the number of voxels and G and S denote the ground truth and segmentation mask, respectively. Normalized Surface Distance (NSD) was used to measure the boundary accuracy, which is defined by

$$NSD(G, S) = \frac{|\partial G \cap B_{\partial S}^{(\tau)}| + |\partial S \cap B_{\partial G}^{(\tau)}|}{|\partial G| + |\partial S|},$$

where $B_{\partial G}^{(\tau)} = \{x \in R^3 \mid \exists \tilde{x} \in \partial G, \|x - \tilde{x}\| \leq \tau\}$ and $B_{\partial S}^{(\tau)} = \{x \in R^3 \mid \exists \tilde{x} \in \partial S, \|x - \tilde{x}\| \leq \tau\}$ denote the border region of the ground truth and the segmentation surface at tolerance τ , respectively. The tolerance τ is defined by measuring organ segmentation consistency between radiologists and revised by oncologists and surgeons based on their clinical requirements (Supplementary). We also found that our boundary tolerances are consistent with other independent inter-rater annotation variability studies [12], [38]. The running time was used to measure the segmentation efficiency which is defined by the duration T between the algorithm docker container start and end. In order to obtain precise metrics, the algorithms were evaluated one by one, and all the other applications on the workstation were closed during the running time. The resource consumption was measured by GPU memory consumption and CPU utilization. Instead of only capturing the maximum consumption [39], we measured all the cumulative resource consumption during the algorithm running time. In particular, we recorded the GPU memory and GPU utilization every 0.1s and defined two new metrics: Area under GPU memory-time curve (AUC_GPU) and Area under CPU utilization-time curve (AUC_CPU), which are defined by

$$AUC_GPU = \sum_{t=0}^T GPU_t$$

and

$$AUC_CPU = \sum_{t=0}^T CPU_t,$$

where GPU_t and CPU_t denote the GPU memory consumption and GPU utilization at timepoint t . If the algorithms get stuck, the corresponding metrics will be set to the worse values (DSC=0, NSD=0, Time=3600, AUC_GPU=3600*(1024*10-2048)=29491200, AUC_CPU=3600*100=360000). Among the 48 evaluated algorithms, only two algorithms got stuck during model inference because they did not optimize the model efficiency (e.g., upsampling the 3D CT scans to high resolution and directly loading the whole volume to memory rather than using a patch-based way).

Ranking scheme

All metrics were used to compute the final ranking. We give GPU memory consumption a tolerance of 2048MB because this kind of GPU is affordable (≤ 100 USD) for most medical centers and personal computers. Thus, the GPU memory consumption at time t was transformed by $G\hat{P}U_t = \max(0, GPU_t - 2048MB)$. We assigned a half weight to CPU and GPU metrics, which can achieve a balance between segmentation accuracy and efficiency&resource consumption in the ranking scheme. The internal validation set contains 200 cases. The corresponding ranking scheme includes three steps:

- Step 1. Computing the five metrics (DSC, NSD, Running Time, AUC_GPU, and AUC_CPU) for each case.
- Step 2. Ranking the 48 algorithms for each case and each metric. Thus, each algorithm will have 200×5 rankings.
- Step 3. Computing the overall rank for each algorithm by averaging all the rankings (AUC_GPU and AUC_CPU metrics are weighted by 0.5).

Ranking stability and statistical analysis

The challenge ranking should be independent of the specific datasets because it reflects the task performances of algorithms. We applied bootstrapping and computed Kendall's τ [40] to quantitatively analyze the variability of our ranking scheme. Specifically, we first extracted 1000 bootstrap samples from the international validation set and computed the ranks again for each bootstrap sample. Then, the ranking agreement was quantified by Kendall's τ . Kendall's τ computes the number of pairwise concordances and discordances between ranking lists. Its value ranges $[-1, 1]$ where -1 and 1 denote inverted and identical order, respectively. A stable ranking scheme should have a high Kendall's τ value that is close to 1.

To compare the performance of different algorithms, we performed Wilcoxon signed rank test because it is a paired comparison. Results were considered statistically significant if the p -value is less than 0.05.

The above analysis was performed with ChallengeR [41], Python 3 [42], Numpy [43], Pandas [44], Scipy [45], Py-Torch [46], and matplotlib [47].

Data availability

The full development set has been released on the challenge website <https://flare22.grand-challenge.org/> and participants can access the data by registering on the Grand Challenge website. The hidden internal validation set and external validation sets will be publicly available on the challenge website after peer review.

Code availability

The code, method description, and docker containers of the top teams are available at <https://flare22.grand-challenge.org/awards/>. The evaluation code is available at <https://github.com/JunMa11/FLARE>.

REFERENCES

- [1] R. L. Siegel, K. D. Miller, H. E. Fuchs, and A. Jemal, "Cancer statistics, 2022," *CA: A Cancer Journal for Clinicians*, vol. 72, no. 1, pp. 7–33, 2022.
- [2] A. A. Perez, V. Noe-Kim, M. G. Lubner, P. M. Graffy, J. W. Garrett, D. C. Elton, R. M. Summers, and P. J. Pickhardt, "Deep learning ct-based quantitative visualization tool for liver volume estimation: defining normal and hepatomegaly," *Radiology*, p. 210531, 2021.
- [3] G. E. Humpire-Mamani, J. Bukala, E. T. Scholten, M. Prokop, B. van Ginneken, and C. Jacobs, "Fully automatic volume measurement of the spleen at ct using deep learning," *Radiology: Artificial Intelligence*, vol. 2, no. 4, p. e190102, 2020.
- [4] J. De Fauw, J. R. Ledsam, B. Romera-Paredes, S. Nikolov, N. Tomasev, S. Blackwell, H. Askham, X. Glorot, B. O'Donoghue, D. Visentin *et al.*, "Clinically applicable deep learning for diagnosis and referral in retinal disease," *Nature Medicine*, vol. 24, no. 9, pp. 1342–1350, 2018.
- [5] D. Ouyang, B. He, A. Ghorbani, N. Yuan, J. Ebinger, C. P. Langlotz, P. A. Heidenreich, R. A. Harrington, D. H. Liang, E. A. Ashley *et al.*, "Video-based ai for beat-to-beat assessment of cardiac function," *Nature*, vol. 580, no. 7802, pp. 252–256, 2020.
- [6] R. Arnaout, L. Curran, Y. Zhao, J. C. Levine, E. Chinn, and A. J. Moon-Grady, "An ensemble of neural networks provides expert-level prenatal detection of complex congenital heart disease," *Nature Medicine*, vol. 27, no. 5, pp. 882–891, 2021.
- [7] C. McIntosh, L. Conroy, M. C. Tjong, T. Craig, A. Bayley, C. Catton, M. Gospodarowicz, J. Helou, N. Isfahanian, V. Kong *et al.*, "Clinical integration of machine learning for curative-intent radiation treatment of patients with prostate cancer," *Nature Medicine*, vol. 27, no. 6, pp. 999–1005, 2021.
- [8] X. Li, H. Chen, X. Qi, Q. Dou, C.-W. Fu, and P.-A. Heng, "H-denseunet: hybrid densely connected unet for liver and tumor segmentation from ct volumes," *IEEE Transactions on Medical Imaging*, vol. 37, no. 12, pp. 2663–2674, 2018.
- [9] N. Heller, F. Isensee, K. H. Maier-Hein, X. Hou, C. Xie, F. Li *et al.*, "The state of the art in kidney and kidney tumor segmentation in contrast-enhanced ct imaging: Results of the kits19 challenge," *Medical Image Analysis*, vol. 67, p. 101821, 2021.
- [10] H. R. Roth, L. Lu, A. Farag, H.-C. Shin, J. Liu, E. B. Turkbey, and R. M. Summers, "Deeporgan: Multi-level deep convolutional networks for automated pancreas segmentation," in *International Conference on Medical Image Computing and Computer-assisted Intervention*, 2015, pp. 556–564.
- [11] J. C. Caicedo, A. Goodman, K. W. Karhohs, B. A. Cimini, J. Ackerman, M. Haghighi, C. Heng, T. Becker, M. Doan, C. McQuin *et al.*, "Nucleus segmentation across imaging experiments: the 2018 data science bowl," *Nature Methods*, vol. 16, no. 12, pp. 1247–1253, 2019.
- [12] M. Antonelli, A. Reinke, S. Bakas, K. Farahani, A. Kopp-Schneider, B. A. Landman, G. Litjens, B. Menze, O. Ronneberger, R. M. Summers *et al.*, "The medical segmentation decathlon," *Nature Communications*, vol. 13, no. 1, pp. 1–13, 2022.
- [13] W. Bulten, K. Kartasalo, P.-H. C. Chen, P. Ström, H. Pinckaers, K. Nagpal, Y. Cai, D. F. Steiner, H. van Boven, R. Vink *et al.*, "Artificial intelligence for diagnosis and gleason grading of prostate cancer: the panda challenge," *Nature Medicine*, vol. 28, no. 1, pp. 154–163, 2022.
- [14] T. Heimann, B. Van Ginneken, M. A. Styner, Y. Arzhaeva, V. Aurich, C. Bauer, A. Beck, C. Becker, R. Beichel, G. Bekes *et al.*, "Comparison and evaluation of methods for liver segmentation from ct datasets," *IEEE Transactions on Medical Imaging*, vol. 28, no. 8, pp. 1251–1265, 2009.
- [15] P. Bilic, P. Christ, H. B. Li, E. Vorontsov, A. Ben-Cohen, G. Kaissis, A. Szeskin, C. Jacobs, G. E. H. Mamani, G. Chartrand *et al.*, "The liver tumor segmentation benchmark (lits)," *Medical Image Analysis*, p. 102680, 2022.

- [16] N. Heller, S. McSweeney, M. T. Peterson, S. Peterson, J. Rickman, B. Stai, R. Tejpal, M. Oestreich, P. Blake, J. Rosenberg, K. Moore, E. Walczak, Z. Rengel, Z. Edgerton, R. Vasdev, A. Kalapara, N. J. Sathianathan, N. Papanikolopoulos, and C. J. Weight, "An international challenge to use artificial intelligence to define the state-of-the-art in kidney and kidney tumor segmentation in ct imaging," *Journal of Clinical Oncology*, vol. 38, no. 6, pp. 626–626, 2020.
- [17] L. Maier-Hein, A. Reinke, M. Kozubek, A. L. Martel, T. Arbel, M. Eisenmann, A. Hanbury, P. Jannin, H. Müller, S. Onogur *et al.*, "Bias: Transparent reporting of biomedical image analysis challenges," *Medical Image Analysis*, vol. 66, p. 101796, 2020.
- [18] J. Ma, B. Wang, and S. Bharadwa, "Fast and Low-resource Semi-supervised Abdominal Organ Segmentation in CT," Mar. 2022. [Online]. Available: <https://doi.org/10.5281/zenodo.6362374>
- [19] T. Radsch, A. Reinke, V. Weru, M. D. Tizabi, N. Schreck, A. E. Kavur, B. Pekdemir, T. Roß, A. Kopp-Schneider, and L. Maier-Hein, "Labeling instructions matter in biomedical image analysis," *Nature Machine Intelligence*, 2023.
- [20] L. Maier-Hein, M. Eisenmann, A. Reinke, S. Onogur, M. Stankovic, P. Scholz, T. Arbel, H. Bogunovic, A. P. Bradley, A. Carass *et al.*, "Why rankings of biomedical image analysis competitions should be interpreted with care," *Nature Communications*, vol. 9, no. 1, pp. 1–13, 2018.
- [21] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2020.
- [22] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *2016 fourth international conference on 3D vision (3DV)*, 2016, pp. 565–571.
- [23] J. Ma, J. Chen, M. Ng, R. Huang, Y. Li, C. Li, X. Yang, and A. Martel, "Loss odyssey in medical image segmentation," *Medical Image Analysis*, vol. 71, p. 102035, 2021.
- [24] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations*, 2014.
- [25] Z. Huang, H. Wang, J. Ye, J. Niu, C. Tu, Y. Yang, S. Du, Z. Deng, L. Gu, and J. He, "Revisiting nnu-net for iterative pseudo labeling and efficient sliding window inference," in *Fast and Low-Resource Semi-supervised Abdominal Organ Segmentation*, J. Ma and B. Wang, Eds. Cham: Springer Nature Switzerland, 2022, pp. 178–189.
- [26] F. Isensee, P. F. Jäger, S. A. A. Kohl, J. Petersen, and K. H. Maier-Hein, "nnu-net: a self-configuring method for deep learning-based biomedical image segmentation," *Nature Methods*, vol. 18, no. 2, pp. 203–211, 2021.
- [27] E. Wang, Y. Zhao, and Y. Wu, "Cascade dual-decoders network for abdominal organs segmentation," in *Fast and Low-Resource Semi-supervised Abdominal Organ Segmentation*, J. Ma and B. Wang, Eds. Cham: Springer Nature Switzerland, 2022, pp. 202–213.
- [28] L. Maier-Hein, A. Reinke, E. Christodoulou, B. Glocker, P. Godau, F. Isensee, J. Kleesiek, M. Kozubek, M. Reyes, M. A. Riegler *et al.*, "Metrics reloaded: Pitfalls and recommendations for image analysis validation," *arXiv preprint arXiv:2206.01653*, 2022.
- [29] L. Seyyed-Kalantari, H. Zhang, M. B. McDermott, I. Y. Chen, and M. Ghassemi, "Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations," *Nature Medicine*, vol. 27, no. 12, pp. 2176–2182, 2021.
- [30] S. Bakas, M. Reyes, A. Jakab, S. Bauer, M. Rempfler, A. Crimi *et al.*, "Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the brats challenge," *arXiv preprint arXiv:1811.02629*, 2018.
- [31] K. A. Goodman, W. F. Regine, L. A. Dawson, E. Ben-Josef, K. Haustermans, W. R. Bosch, J. Turian, and R. A. Abrams, "Radiation therapy oncology group consensus panel guidelines for the delineation of the clinical target volume in the postoperative treatment of pancreatic head cancer," *International Journal of Radiation Oncology* Biology* Physics*, vol. 83, no. 3, pp. 901–908, 2012.
- [32] F. H. Netter, *Atlas of human anatomy*. Elsevier Health Sciences, 2014.
- [33] J. Ma, Y. Zhang, S. Gu, C. Zhu, C. Ge, Y. Zhang, X. An, C. Wang, Q. Wang, X. Liu, S. Cao, Q. Zhang, S. Liu, Y. Wang, Y. Li, J. He, and X. Yang, "Abdoment-1k: Is abdominal organ segmentation a solved problem?" *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 6695–6714, 2022.
- [34] B. Landman, Z. Xu, J. Igelsias, M. Styner, T. Langerak, and A. Klein, "Miccai multi-atlas labeling beyond the cranial vault-workshop and challenge," 2015.
- [35] X. Luo, W. Liao, J. Xiao, J. Chen, T. Song, X. Zhang, K. Li, D. N. Metaxas, G. Wang, and S. Zhang, "Word: A large scale dataset, benchmark and clinical applicable study for abdominal organ segmentation from ct image," *Medical Image Analysis*, vol. 82, p. 102642, 2022.
- [36] Y. Ji, H. Bai, J. Yang, C. Ge, Y. Zhu, R. Zhang, Z. Li, L. Zhang, W. Ma, X. Wan *et al.*, "Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation," *arXiv preprint arXiv:2206.08023*, 2022.
- [37] J. Wasserthal, H.-C. Breit, M. T. Meyer, M. Pradella, D. Hinck, A. W. Sauter, T. Heye, D. Boll, J. Cyriac, S. Yang, M. Bach, and M. Segeroth, "Totalsegmentator: Robust segmentation of 104 anatomic structures in ct images," *Radiology: Artificial Intelligence*, vol. 0, p. e230024, 2023.
- [38] P. M. Adamson, V. Bhattbhatt, S. Principi, S. Beriwal, L. S. Strain, M. Offe, A. S. Wang, N.-J. Vo, T. Gilat Schmidt, and P. Jordan, "Evaluation of a v-net autosegmentation algorithm for pediatric ct scans: Performance, generalizability, and application to patient-specific ct dosimetry," *Medical Physics*, vol. 49, no. 4, pp. 2342–2354, 2022.
- [39] J. Ma, Y. Zhang, S. Gu, X. An, Z. Wang, C. Ge, C. Wang, F. Zhang, Y. Wang, Y. Xu, S. Gou, F. Thaler, C. Payer, D. Štern, E. G. Henderson, D. M. McSweeney, A. Green, P. Jackson, L. McIntosh, Q.-C. Nguyen, A. Qayyum, P.-H. Conze, Z. Huang, Z. Zhou, D.-P. Fan, H. Xiong, G. Dong, Q. Zhu, J. He, and X. Yang, "Fast and low-gpu-memory abdomen ct organ segmentation: The flare challenge," *Medical Image Analysis*, vol. 82, p. 102616, 2022.
- [40] M. G. Kendall, "A new measure of rank correlation," *Biometrika*, vol. 30, no. 1/2, pp. 81–93, 1938.
- [41] M. Wiesenfarth, A. Reinke, B. A. Landman, M. Eisenmann, L. A. Saiz, M. J. Cardoso *et al.*, "Methods and open-source toolkit for analyzing and visualizing challenge results," *Scientific Reports*, vol. 11, no. 1, pp. 1–15, 2021.
- [42] G. Van Rossum and F. L. Drake, *Python 3 Reference Manual*. Scotts Valley, CA: CreateSpace, 2009.
- [43] C. R. Harris, K. J. Millman, S. J. Van Der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith *et al.*, "Array programming with numpy," *Nature*, vol. 585, no. 7825, pp. 357–362, 2020.
- [44] T. pandas development team, "pandas-dev/pandas: Pandas," 2020. [Online]. Available: <https://doi.org/10.5281/zenodo.3509134>
- [45] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright *et al.*, "Scipy 1.0: fundamental algorithms for scientific computing in python," *Nature methods*, vol. 17, no. 3, pp. 261–272, 2020.
- [46] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems*, 2019, pp. 8024–8035.
- [47] J. D. Hunter, "Matplotlib: A 2d graphics environment," *Computing in Science & Engineering*, vol. 9, no. 3, pp. 90–95, 2007.