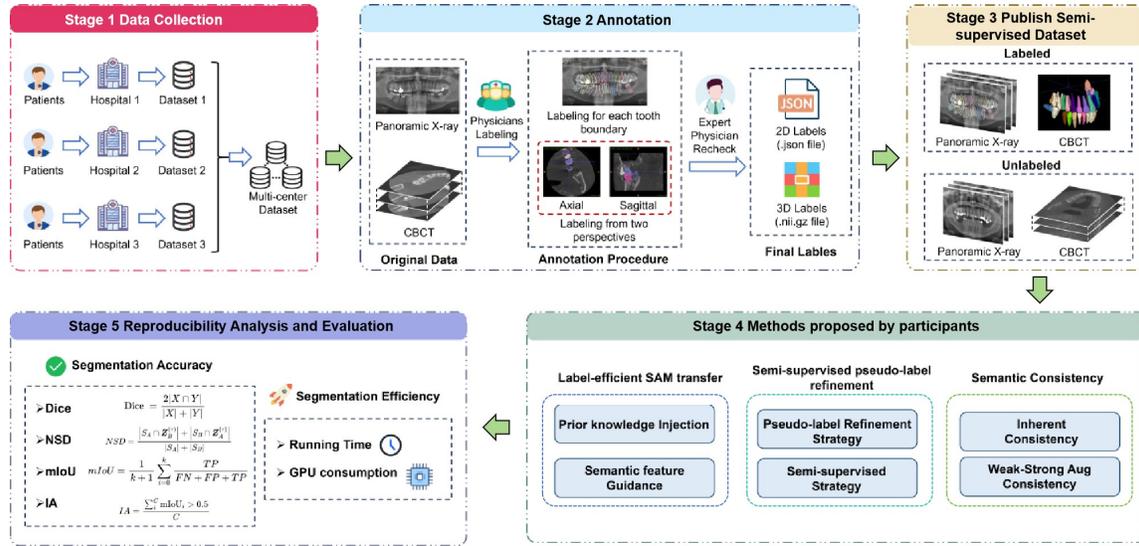


Graphical Abstract



End-to-end workflow of the MICCAI 2024 Semi-supervised Teeth Segmentation (STS) Challenge. The process encompasses five key stages: (1) multi-center data collection to ensure dataset diversity, (2) iterative annotation by clinicians for high-quality ground truth, (3) construction of the semi-supervised dataset with distinct labeled and unlabeled sets, and (4/5) the final summarization and evaluation of submitted participant methods.

Highlights

- A novel dataset for semi-supervised instance-level tooth segmentation.
- First international benchmark of SSL for tooth instance segmentation.
- SSL boosts instance segmentation accuracy by over 60 percentage points compared to a supervised baseline.
- Analysis of top methods including SAM integration and pseudo-labeling.
- Provides insights for label-efficient AI in clinical dental imaging.

MICCAI STS 2024 Challenge: Semi-Supervised Instance-Level Tooth Segmentation in Panoramic X-ray and CBCT Images

Yaqi Wang^{a,1}, Zhi Li^{b,1}, Chengyu Wu^g, Jun Liu^{a,*}, Yifan Zhang^d, Jiaxue Ni^c, Qian Luo^c, Jialuo Chen^c, Hongyuan Zhang^e, Jin Liu^c, Can Hanⁱ, Kaiwen Fu^h, Changkai Ji^j, Xinxu Cai^j, Jing Hao^k, Zhihao Zheng^l, Shi Xu^m, Junqiang Chenⁿ, Qianni Zhang^f, Dahong Qianⁱ, Shuai Wang^{b,*}, Huiyu Zhou^{o,*}

^a*Innovation Center for Electronic Design Automation Technology, Hangzhou Dianzi University, Hangzhou, China*

^b*School of Cyberspace, Hangzhou Dianzi University, Hangzhou, China*

^c*Hangzhou Dianzi University, Hangzhou, China*

^d*Hangzhou Geriatric Stomatology Hospital, Hangzhou Dental Hospital Group, Hangzhou, China*

^e*Shenzhen University, Shenzhen, China*

^f*Queen Mary University of London, London, United Kingdom*

^g*Shandong University, Weihai, China*

^h*Xidian University, Xi'an, China*

ⁱ*Shanghai Jiao Tong University, Shanghai, China*

^j*Harbin Institute of Technology, Harbin, China*

^k*The University of Hong Kong, Hong Kong SAR, China*

^l*Chinese Academy of Sciences, Chengdu, China*

^m*Yunnan Provincial Stomatology Hospital, Kunming, China*

ⁿ*Shanghai MediWorks Precision Instruments Co., Ltd, Shanghai, China*

^o*University of Leicester, Leicester, United Kingdom*

Abstract

Orthopantomogram (OPGs) and Cone-Beam Computed Tomography (CBCT) are vital for dentistry, but creating large datasets for automated tooth segmentation is hindered by the labor-intensive process of manual instance-level annotation. This research aimed to benchmark and advance semi-supervised learning (SSL) as a solution for this data scarcity problem. We organized the 2nd Semi-supervised Teeth Segmentation (STS 2024) Challenge at MICCAI 2024. We provided a large-scale

*Corresponding authors: ljun77@hdu.edu.cn (Jun Liu); shuaiwang.tai@gmail.com (Shuai Wang), hz143@leicester.ac.uk (Huiyu Zhou)

¹These authors contributed equally to this work.

dataset comprising over 90,000 2D images and 3D axial slices, which includes 2,380 OPG images and 330 CBCT scans, all featuring detailed instance-level FDI annotations on part of the data. The challenge attracted 114 (OPG) and 106 (CBCT) registered teams. To ensure algorithmic excellence and full transparency, we rigorously evaluated the valid, open-source submissions from the top 10 (OPG) and top 5 (CBCT) teams, respectively. All successful submissions were deep learning-based SSL methods. The winning semi-supervised models demonstrated impressive performance gains over a fully-supervised nnU-Net baseline trained only on the labeled data. For the 2D OPG track, the top method improved the Instance Affinity (IA) score by over 44 percentage points. For the 3D CBCT track, the winning approach boosted the Instance Dice score by 61 percentage points. This challenge confirms the substantial benefit of SSL for complex, instance-level medical image segmentation tasks where labeled data is scarce. The most effective approaches consistently leveraged hybrid semi-supervised frameworks that combined knowledge from foundational models like SAM with multi-stage, coarse-to-fine refinement pipelines. Both the challenge dataset and the participants' submitted code have been made publicly available on GitHub (<https://github.com/ricoleehduu/STS-Challenge-2024>), ensuring transparency and reproducibility.

Keywords: Tooth Segmentation, Semi-supervised Learning, CBCT, OPG

1. Introduction

Oral health is integral to overall well-being, with dental structures profoundly impacting an individual's quality of life [1]. Modern dentistry relies heavily on advanced imaging modalities such as 2D Orthopantomogram (OPGs) and 3D Cone-Beam Computed Tomography (CBCT) [2]. OPGs offer a broad view of the dentition and surrounding structures, while CBCT provides high-resolution 3D data for detailed anatomical evaluation, crucial for detecting caries, impacted or supernumerary teeth, and for precise diagnosis and treatment planning in orthodontics and implantology [3, 4]. Accurate tooth segmentation from these images is fundamental for many computer-aided systems, enhancing clinical decision making, particularly when performed at the instance level [5]. Therefore, robust automated instance-level tooth segmentation algorithms are essential to advance the precision, efficiency, and personalization of dental care [6], a process visualized in Fig. 1.

Despite its clinical importance, the development of high-performance automated tooth segmentation models is critically hampered by data scarcity. Manual annotation, especially for detailed instance-level identification, is extremely time-consuming, laborious, and requires expert knowledge [7]. This results in a critical shortage of

large-scale, high-quality labeled datasets essential for training conventional fully supervised deep learning models. Consequently, while fully supervised methods show promise, their performance is often limited by the lack of diverse and comprehensively annotated data, especially for complex cases.

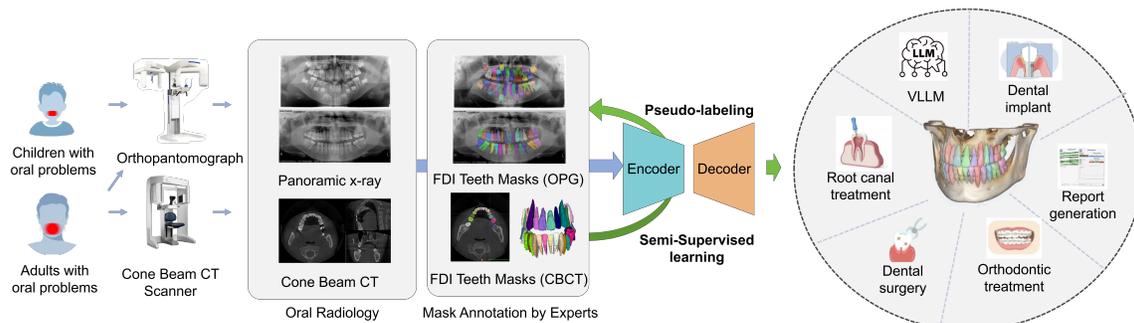


Figure 1: Overview of the Semi-Supervised Learning (SSL) framework for dental instance segmentation and its clinical utility. The workflow proceeds from (Left) multi-modal data acquisition (OPG and CBCT) covering both pediatric and adult populations, to (Middle) the core semi-supervised training paradigm where an Encoder-Decoder network leverages expert annotations and iteratively refines performance via pseudo-labeling, and finally to (Right) diverse downstream clinical applications, such as orthodontic treatment planning, root canal therapy, and automated report generation, which rely on precise tooth instance masks.

Semi-supervised learning (SSL) has emerged as a promising paradigm to address limited labeled data in medical image analysis, including dental imaging [8]. SSL methods leverage a small labeled dataset alongside a larger pool of unlabeled images, reducing reliance on extensive manual annotation and potentially improving model generalization and robustness. However, SSL presents its own challenges. A core difficulty is generating and refining high-quality pseudo-labels for unlabeled data, as inaccuracies can propagate errors and degrade performance. Models may also suffer from confirmation bias. Thus, effective SSL requires careful design of consistency regularization, uncertainty estimation, and pseudo-label selection mechanisms to ensure genuine learning from unlabeled data without over-reliance on the initial labeled set [9].

Tooth segmentation becomes considerably more complex when progressing from semantic delineation (tooth vs. background) to instance segmentation. Instance-level segmentation demands not only precise boundary delineation but also unique identification of each tooth, often with its FDI number, effectively treating each tooth instance as a distinct class. This complexity is amplified by variable dental anatomy, such as morphology, arrangement, crowding, and impaction, across diverse age groups and pathologies. Close tooth proximity and overlap, especially in 2D

OPGs, hinder accurate instance separation. Image quality issues in OPG and CBCT, like distortions, artifacts, and varying contrast, add further difficulty [10]. Moreover, integrating robust FDI numbering within an SSL framework with sparse ground truth for such identification is a significant hurdle, compounded by data imbalances and the need for generalization across varied clinical settings and acquisition protocols.

Table 1: Summary of labeled and unlabeled data statistics for the 2D dental panoramic X-ray and 3D Cone-Beam Computed Tomography (CBCT) datasets, including patient/scans counts, image resolution, and anatomical annotations.

Dataset	Metric	Labeled	Unlabeled
OPG(2D)	Child Samples	30	887
	Adult Samples	70	1484
	Number of Patients	100	2323
	Resolutions (Child)(Pixel)	2000×942	2000×942
	Resolutions (Adult)(Pixel)	1991×1127	2000×942 (887) / 1991×1127 (1463)
	Number of Teeth	≈ 2700	≈ 65044
CBCT(3D)	Adult Samples	100	300
	Number of Patients	100	300
	Number of Slices	29240	60000
	Resolutions(Voxel)	$266 \times 266 / 512 \times 512$	266×266
	In-plane Resolution (mm)	$0.3 \times 0.3 / 0.25 \times 0.25$	0.3×0.3
	Slice Thickness (mm)	$0.3 / 0.25$	0.3
	Number of Teeth	≈ 3000	≈ 9000

Despite the recognized need for automated dental image analysis, developing and benchmarking advanced algorithms, especially for instance-level tasks, is hindered by limited public datasets. Although there are several dental datasets [11], many have limitations for robust semi-supervised instance-level systems (Table 2, 3). For example, some focus on adult populations or lack granular FDI-numbered annotations crucial for clinical applications. Critically, most are not designed for SSL, often missing large, matched, unlabeled image sets. Recognizing these gaps, we initiated the Semi-supervised Teeth Segmentation (STS) challenge series. Although our inaugural STS 2023 challenge established a baseline for semantic segmentation [12], the more clinically crucial and complex task of instance-level segmentation remained an open problem, motivating the current work.

Building on this foundation, the 2nd Semi-supervised Teeth Segmentation (STS 2024) Challenge was organized as an official satellite event of the 27th International Conference on Medical Image Computing and Computer-Assisted Intervention (MIC-CAI 2024). Its primary aim was to stimulate the development and rigorous evalua-

tion of novel SSL algorithms for instance-level tooth segmentation in both 2D OPGs and 3D CBCTs, creating a standardized benchmark for this challenging task. The challenge featured two tasks (OPG and CBCT), requiring semi-supervised instance segmentation and identification of up to 32 permanent and 20 deciduous teeth. Participants received datasets with a small fraction of instance-level annotations and a majority of unlabeled data as described in Table 1, with evaluations covering both accuracy and efficiency. The challenge dataset, including the labeled training set, validation set, and the large-scale unlabeled set, is now publicly available for direct download in our GitHub repository (<https://github.com/ricoleehduu/STS-Challenge-2024>). The complete workflow of the STS 2024 Challenge is shown in Fig. 2.

In general, the main contributions are summarized as follows:

- A novel public dataset for semi-supervised instance-level tooth segmentation, the first large-scale, multi-modal resource including OPG and CBCT, specifically curated for semi-supervised instance-level tooth segmentation. The dataset features detailed FDI annotations for a small labeled set alongside extensive unlabeled data, covering diverse patient demographics and pathologies.
- We established the first open international benchmark for this specific task, employing a rigorous evaluation framework with Dockerized submissions on a hidden test set to ensure fair and reproducible comparisons of state-of-the-art methods.
- Provide a detailed analysis and summary of the various SSL strategies employed by the top-performing participants. This includes an overview of popular architectures, effective SSL techniques, and data handling approaches, offering valuable insights to guide future research in label-efficient AI for dentistry.

2. Related Works

2.1. Datasets and Benchmarks in Dental Image Segmentation

Artificial intelligence (AI) is increasingly vital in stomatology, where AI-driven segmentation of teeth and maxillofacial structures from Orthopantomogram (OPGs) and Cone-Beam Computed Tomography (CBCT) is a critical enabler for applications like lesion visualization, orthodontic design, and surgical planning [26]. This has led to the creation of numerous public datasets to fuel research in this area (see Tables 2 and 3). However, a closer examination reveals critical gaps, particularly for developing robust, label-efficient, instance-level segmentation models.

Table 2: Comparison of the STS 2024 dataset with other public 2D OPG datasets. The STS 2024 dataset provides a unique resource for semi-supervised learning, featuring a large unlabeled set, multi-center data, and detailed instance-level annotations across a wide age range. "A / woA" denotes data with/without annotations.

Dataset	Year	X-Ray Slices (A / woA)	Num. of Teeth	Annotation Type	Annotators	Centers	Resolution (Pixels)
Dental Panoramic X-Rays [13]	2017	116 (116 / 0)	$\approx 3,480$	Tooth Instances	2	1	3100×1300
UFBA-UESC [14]	2019	1,500 (543 / 957)	$\approx 45,000$	Instance, Numbering	1	1	1991×1127
Tufts Dental [15]	2021	1,000 (1000 / 0)	$\approx 27,000$	Semantic, Caries Attn.	2	1	1615×840
Panoramic Radiography [16]	2021	598 (0 / 598)	$\approx 17,940$	Lesion Segmentation	1	1	2041×1024
DENTEX [17]	2023	2,332 (1005 / 1227)	$\approx 69,960$	Tooth Bounding Box	3	1	2098×970
vzrad2 [18]	2024	8,188 (8188 / 0)	$\approx 245,640$	Pathology, Treatment	1	1	640×640
STS2023 [12]	2023	6,500 (6500 / 0)	$\approx 195,000$	Semantic	30	2	640×640
STS2024 (Ours)	2024	2,380 (30 / 2350)	$\approx 71,400$	Instance, Numbering	30	2	2000×942 1991×1127

Many existing datasets, while valuable, have significant limitations for our task. For instance, some primarily offer semantic-level annotations, which cannot distinguish between individual teeth. The Tufts Dental [15] dataset and our own previous STS2023 [12] dataset fall into this category, providing strong baselines for tooth-background separation but not for instance identification. Other datasets provide instance-level information, but with limitations in scope or annotation type. For example, DENTEX [17] provides only bounding boxes rather than fine-grained masks, and datasets like UFBA-UESC [14] and CTooth+ [20] focus on specific populations, such as excluding pediatric teeth or having a narrow age range, which limits model generalizability to diverse clinical scenarios.

While theoretically any fully annotated dataset can be adapted for semi-supervised learning (SSL) by withholding labels, the specific challenge of dental instance segmentation demands a dataset structure that mirrors real-world clinical constraints. The meticulous, pixel-level instance annotation required for tasks like root canal therapy and maxillofacial surgery creates a severe annotation bottleneck. Our STS 2024 dataset addresses this by providing a curated benchmark specifically designed to test SSL algorithms on a difficult multi-class instance segmentation task. It combines high-quality, expert-verified instance masks for a small labeled set with a massive, clinically representative corpus of unlabeled data, directly simulating the scenario

Table 3: Comparison of STS2024 dataset and other public 3D dental CBCT datasets. The STS2024 dataset provides broader age coverage, more detailed annotations, and multi-center diversity, making it a valuable resource for clinical dental AI research.

Dataset	Years	Age Range	Patients	Volumes (A / woA)	Slices	Num of Teeth	Annotator	Volumes (pixels)
Clinical dental CBCT [19]	2020	10-49	25	25 (25 / 0)	9400	≈ 770	1	110 × 145 × 280
CTooth+[20]	2022	10–15	22	168 (22 / 146)	31380	≈ 5040	15	266 × 266 × 266
Mandibular Canal CBCT [21]	2022	10-100	347	347 (347 / 0)	88832	≈ 17220	-	178 × 423 × 463
Multi-modal dataset [22]	2023	18-89	169	188 (188 / 188)	16203	≈ 5401	13	640 × 640 × 200
STS2023 [12]	2023	7–70	584	584 (84 / 500)	88500	≈ 17520	30	640 × 640 × 399
ToothFairy[23]	2022	12-100	443	443 (443 / 0)	64000	– (IAC only)	5	169 × 342 × 370
ToothFairy2[24]	2023	16-100	530	530 (480/50)	80000	≈ 71400	5	170 × 272 × 345 298 × 512 × 512
ToothFairy3[25]	2024	16-100	532	532 (532/0)	80000	≈ 17024	5	170 × 272 × 345 298 × 512 × 512
STS2024 (Ours)	2024	7–70	330	330 (30 / 300)	69960	≈ 9900	30	266 × 266 × 200 512 × 512 × 332

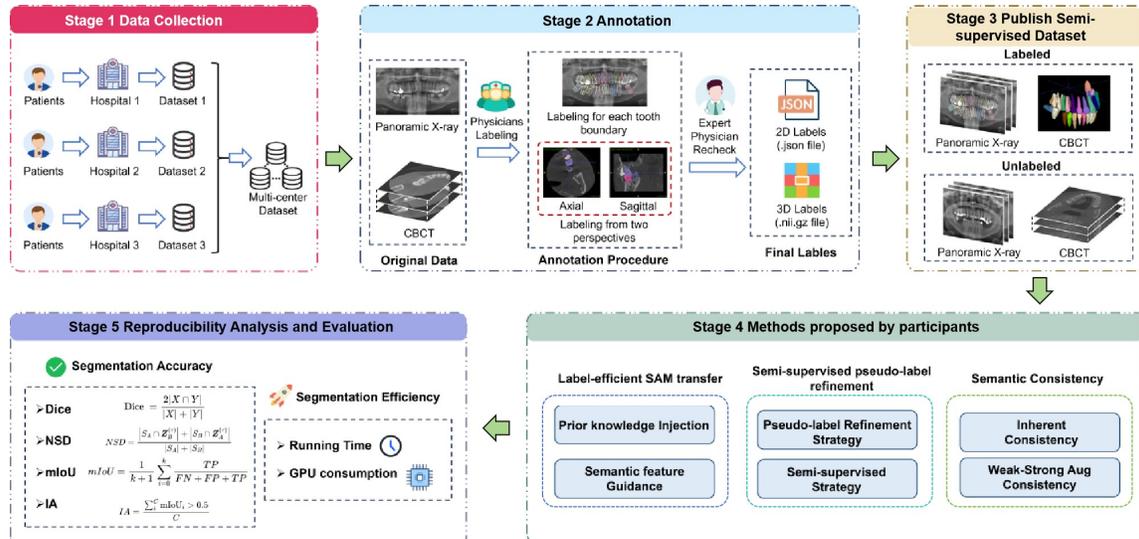


Figure 2: End-to-end workflow of the MICCAI 2024 Semi-supervised Teeth Segmentation (STS) Challenge. The process encompasses five key stages: (1) multi-center data collection to ensure dataset diversity, (2) iterative annotation by clinicians for high-quality ground truth, (3) construction of the semi-supervised dataset with distinct labeled and unlabeled sets, and (4/5) the final summarization and evaluation of submitted participant methods.

where leveraging abundant raw data is the only viable path to high performance.

To directly address these multifaceted limitations, we initiated the Semi-supervised Teeth Segmentation (STS) challenge series. Although our inaugural STS 2023 challenge [12] established a robust baseline for SSL-based semantic segmentation, the ability to distinguish and number individual teeth, a prerequisite for most advanced clinical applications, remained an unaddressed challenge. To bridge this critical gap, the STS 2024 Challenge was designed to elevate the task to semi-supervised instance-level segmentation. The challenge introduces a unique, multimodal dataset with detailed, instance-level FDI annotations complemented by a vast pool of unlabeled OPG and CBCT data. Spanning a wide range of ages (7–70 years) and pathologies, this competitive framework is the first of its kind designed to catalyze and rigorously evaluate SSL algorithms for precise, instance-aware segmentation, with the aim of advancing intelligent and efficient dental care.

2.2. Methodologies for Tooth Segmentation and Semi-Supervised Learning

Automated tooth segmentation has progressed from traditional image processing to deep learning (DL). In fully supervised settings, U-Net [27] and its variants like V-Net [28] and nnU-Net [29] are standard for semantic segmentation in dental imaging [30, 31]. For instance-level segmentation, which distinguishes individual teeth, common approaches adapt computer vision methods like Mask R-CNN [32] or use post-processing techniques to separate instances after semantic segmentation. However, the performance of these fully supervised methods is fundamentally tethered to the availability of large, meticulously annotated datasets, a condition rarely met in specialized medical domains like dentistry.

Semi-supervised learning (SSL) mitigates this dependency on labeled data. One prominent SSL paradigm is pseudo-labeling, also known as self-training. In this approach, a model first trains on a small labeled set and then generates predictions, or pseudo-labels, for unlabeled data. These new labels are then used to augment the training set, allowing the model to learn from the unlabeled images [33]. A critical challenge here is managing the quality of these pseudo-labels to prevent error propagation. This is often addressed through techniques like confidence thresholding and iterative refinement [34]. However, in dense dental anatomies, generating high-confidence pseudo-labels for heavily overlapping or small, developing teeth remains a significant challenge. Another major SSL family is consistency regularization. The core principle is that a model’s prediction for an unlabeled input should remain consistent even when the input is perturbed, for example, through data augmentation or internal model changes like dropout [35]. The Mean Teacher model is a classic example of this, where a "student" network is trained to produce predictions consistent

with a "teacher" network, which is an exponential moving average of the student's own weights [35]. A key challenge for consistency-based methods in dental imaging is defining realistic augmentations for OPG/CBCT data without altering the fine anatomical boundaries critical for segmentation.

Many state-of-the-art SSL methods combine these two strategies. They often use teacher-student frameworks to guide learning on unlabeled data via robust pseudo-label generation [36, 37]. Other techniques like multi-stage training and uncertainty-aware learning are also employed to enhance performance [38]. While SSL has been successful in various medical segmentation tasks [39, 40], its application to instance-level dental segmentation was relatively limited before initiatives like the STS challenges [38].

Emerging DL trends are also shaping this field. Foundation models, such as the Segment Anything Model (SAM) [41], provide powerful pre-trained capabilities that can be used for few-shot segmentation or to generate high-quality initial pseudo-labels, potentially overcoming the difficulty of separating overlapping teeth where traditional models fail [42]. Another effective strategy is self-supervised pre-training on large unlabeled datasets before fine-tuning with SSL [37]. Finally, automated frameworks like nnU-Net [29] serve as strong baselines that can be adapted for SSL. The STS 2024 challenge, therefore, provides a timely and crucial benchmark to evaluate how these diverse SSL strategies and emerging techniques contend with the real-world complexities of instance-level tooth segmentation.

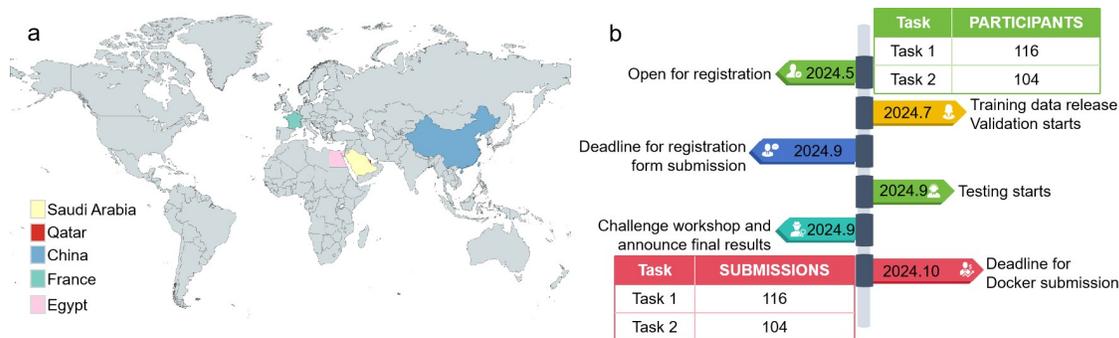


Figure 3: Overview of the STS 2024 Challenge participation and timeline. (a) A world map illustrating the geographical distribution of registered participants. (b) A detailed timeline of the challenge schedule, from the training phase start to the final announcement of results.

3. Challenge Description

3.1. Dataset information and annotation details

The STS (2nd Semi-Supervised Tooth Segmentation) Challenge is one of the official challenges of MICCAI 2024, which aims to advance SSL-based tooth image segmentation, focusing on algorithmic performance in multi-instance, multi-class scenarios. The challenge required participants to leverage a small labeled dataset and a large unlabeled dataset to achieve accurate, instance-level tooth segmentation and classification. The challenge consists of two tasks: Task 1 is based on 2D panoramic radiographs (OPG) and Task 2 is based on 3D cone beam CT (CBCT), where up to 32 permanent teeth and 20 deciduous teeth need to be accurately recognized, reflecting the full spectrum of possible dentition. The challenge was hosted on the Codabench platform and attracted global participation, as illustrated by the participant distribution map and competition timeline in Fig. 3. The schedule was divided into distinct phases, including training, validation, and a final testing phase on a hidden dataset to ensure a fair and rigorous evaluation of all submitted methods.

The datasets used in this challenge were provided by Hangzhou Stomatological Hospital and Hangzhou Qiantang Stomatological Hospital, and were from real clinical environments, covering a wide range of pre-treatment images such as missing teeth and orthodontic appliances. All images were acquired by radiologists or dentists with more than five years of experience and labeled by a team of 30 dentists following strict Standard Operating Procedures (SOPs). To ensure consistency, a pilot annotation phase was conducted to unify boundary definitions across the team. A two-stage annotation protocol was employed: each scan was initially annotated by one of ten junior dentists and subsequently reviewed and corrected by one of three senior dentists with over 10 years of experience. To validate the reliability of the ground truth, we assessed the inter-observer agreement on a randomly selected subset of the data. Any remaining discrepancies were resolved through panel consensus. The dataset has been approved by the Medical Ethics Committee and is licensed under the CC BY-NC-ND license for scientific use only, and commercialization and secondary distribution are prohibited. The data for each task consisted of a small labeled training set and a large unlabeled set used to support the development and evaluation of semi-supervised learning models.

3.1.1. 2D-XRay Dataset

For Task 1 (2D Panoramic X-ray Dataset), images were categorized into adult and child subsets based on patient age and tooth morphology. During preprocessing, the original DICOM images were converted to PNG format. Dental experts performed instance-level annotations using tools such as EISeg and LabelMe, with a standard

image resolution of 640×320 pixels. The annotations, including FDI numbering for each tooth instance, were saved as JSON files. The dataset was divided into training, validation, and test sets, maintaining a realistic clinical ratio of adult to pediatric patients. The training set for Task 1 consisted of 2380 OPGs (30 fully annotated for instance-level segmentation, 2350 unlabeled), with a validation set of 20 OPGs and a test set of 50 OPGs.

3.1.2. 3D-CBCT Dataset

For Task 2 (3D-CBCT dataset), each 3D volume was meticulously annotated layer by layer by dental experts using ITK-SNAP software. All identifiable teeth (primarily permanent teeth in this task) within each volume were segmented at the instance level and assigned corresponding FDI tooth numbers. The annotation results were saved as nii.gz files. The training set for Task 2 includes 330 CBCT scans (30 of which are fully annotated at the instance level, and 300 are unannotated), the validation set includes 20 CBCT scans, and the test set includes 50 CBCT scans. Among the unlabelled samples, 62.79 % had artefacts, 54.47% had fillings, 88.15% had missing teeth, 92.72% had implants, 94.59% had decayed teeth/roots, 97.92% had deciduous teeth, and 98.34% had dental braces. Even in the unlabelled training set samples, 35.00% still contained artefacts.

3.2. Participants and challenge phases

The STS2024 Challenge is aiming to promote the research and practical application of semi-supervised learning methods in medical image analysis. The challenge is conducted in phases: registration for the competition opens on May 10, 2024, training data is released, and the online validation phase is launched on July 15, test data opens on September 23, and participants are required to submit their final algorithms based on the Docker package before September 26th. The final evaluation results are officially released through the STS2024 Challenge Workshop during the MICCAI 2024 conference, which was held October 6-10, 2024, in Marrakech, Morocco.

The Challenge attracted research teams from several countries and regions around the world, with the largest number of participants from China. According to the participation statistics, Task 1 (2D tooth segmentation task) was participated in by 116 teams with a cumulative number of 346 submissions; Task 2 (3D CBCT tooth segmentation task) was participated in by 104 teams with a cumulative number of 158 submissions during the preliminary phase. However, to uphold the highest standards of reproducibility, the final leaderboard focuses on the 5 teams that successfully submitted valid Docker containers and open-sourced their code for the hidden test set evaluation. These data reflect that researchers around the world have paid great

attention to tooth image segmentation technology and extensively explored semi-supervised learning methods in real applications. To recognize outstanding contributions, cash prizes of \$500 were awarded to the top-performing team in each track, with additional souvenirs for other top teams presenting in person.

These data reflect the significant global attention on tooth image segmentation technology and the extensive exploration of semi-supervised learning methods. To foster collaboration and reproducibility, the solutions from the top 10 teams in the 2D track and top 5 in the 3D track are publicly available on our GitHub repository, accompanied by detailed result comparisons. Furthermore, the first authors of the leading teams were invited to co-author this summary paper. Throughout the challenge, participants submitted Dockerized algorithms for evaluation, with a multi-phase validation process allowing them to test and refine their methods multiple times before the final submission, ensuring a fair and robust assessment.

3.3. Clinical Utility of Segmentation

The clinical relevance of this challenge arises from its focus on accurate instance-level segmentation of teeth, which has far-reaching clinical implications for various dental specialties. In addition to general identification of tooth regions, the precise localization of each tooth using FDI numbers allows for finer and more accurate diagnosis and treatment.

In orthodontics, precise segmentation of tooth instances [43] is an important prerequisite for diagnosing malocclusion, planning aligner placement paths, and evaluating treatment outcomes. By accurately modeling tooth morphology and spatial relationships, the surgeon can more effectively formulate movement paths and mechanical strategies to enhance the controllability and effectiveness of treatment.

In terms of dental implant surgical planning [44], instance segmentation can clearly distinguish the target teeth, neighboring tooth roots, and key anatomical structures within the jawbone (e.g., the inferior alveolar nerve), which in turn aids in assessing the bone volume and risk in the implant area. This plays an important role in improving the accuracy of implant position and reducing intraoperative complications. Additionally, in restorative dentistry and prosthetics, obtaining the precise boundaries of each tooth is the basis for digitally customizing the design of crowns, bridges, and removable prostheses, helping to achieve better occlusal fit and aesthetics.

More broadly, instance-level segmentation [45] with AI can also help clinicians visualize lesion progression (e.g., caries, apical periodontitis, etc.), quantitatively assess the evolution of dental disease over time, and even perform 3D visualization simulations prior to complex maxillofacial surgeries. Together, these capabilities

drive dental imaging towards intelligent diagnosis and personalized treatment.

3.4. Performance Evaluation

The competition requires participants to submit segmentation masks generated on the original test images (in json file format for 2D tasks and .nii.gz for 3D tasks), which are evaluated using a variety of performance and efficiency metrics, including image-level versus instance-level Dice Similarity Coefficients (DSCs), Normalized Surface Distance (NSDs), and Recognition Accuracy (IAs). Algorithm runtime (no more than 60 seconds per case) and GPU memory consumption (based on memory-time curve area) will also be examined.

The pixel-level Dice coefficient is a set similarity measure function that is used to evaluate the degree of similarity between two sets, and its formula is defined as follows:

$$\text{Dice}_{\text{image}} = \frac{2 * |A \cap B|}{|A| + |B|} \quad (1)$$

where A denotes the mask of the proposed model prediction and B denotes the mask of Ground Truth (GT) labeling.

The pixel-level Normalized Surface Dice (NSD) measures the overall segmentation quality by quantifying the degree of surface overlap between the predicted and ground-truth boundaries. It is defined as follows:

$$\text{NSD}_{\text{image}} = \frac{\text{overlap}_{\text{GT}} + \text{overlap}_{\text{pred}}}{\text{total_area}_{\text{GT}} + \text{total_area}_{\text{pred}}}, \quad (2)$$

In this equation, $\text{overlap}_{\text{GT}}$ denotes the surface area of the ground-truth boundary that lies within the specified tolerance distance from the predicted boundary, and $\text{overlap}_{\text{pred}}$ represents the surface area of the predicted boundary that lies within the same tolerance distance from the ground-truth boundary. $\text{total_area}_{\text{GT}}$ and $\text{total_area}_{\text{pred}}$ refer to the total surface areas (in pixels) of the ground-truth and predicted boundaries, respectively. In this study, the tolerance distance was fixed at 2 mm, which defines the acceptable spatial deviation between the predicted and ground-truth surfaces. Points on the two surfaces are considered overlapping if their Euclidean distance is less than or equal to 2 mm. A higher NSD value indicates a greater degree of surface consistency and thus better segmentation accuracy at the boundary level.

The instance-level Normalized Surface Dice (NSD) evaluates the segmentation performance across multiple individual objects or instances and is defined as follows:

$$\text{NSD}_{\text{instance}} = \frac{1}{N} \sum_{i=1}^N \text{NSD}_i, \quad (3)$$

$$\text{where } \text{NSD}_i = \frac{\text{overlap}_{\text{GT}_i} + \text{overlap}_{\text{pred}_i}}{\text{total_area}_{\text{GT}_i} + \text{total_area}_{\text{pred}_i}},$$

Here, N denotes the total number of instances in the dataset. For each instance i , $\text{overlap}_{\text{GT}_i}$ and $\text{overlap}_{\text{pred}_i}$ represent the overlapping areas between the predicted and ground-truth boundaries of that specific instance within the 2 mm tolerance range, while $\text{total_area}_{\text{GT}_i}$ and $\text{total_area}_{\text{pred}_i}$ denote their respective total surface areas. The instance-level NSD reflects the average boundary accuracy across all segmented instances, ensuring that both global and object-wise geometric consistency are comprehensively evaluated.

To specifically evaluate the critical task of instance identification and separation, we introduced the Instance Affinity (IA) metric. This metric is defined as the fraction of ground truth tooth instances that are correctly detected, where a detection is considered correct if the Intersection over Union (IoU) between the predicted instance and the ground truth instance is greater than or equal to 0.5.

$$\text{IA} = \frac{\sum_{i=1}^{N_{\text{GT}}} \mathbb{I}(\max_j(\text{IoU}(G_i, P_j)) \geq 0.5)}{N_{\text{GT}}} \quad (4)$$

where G_i is the i -th ground truth instance, P_j is a predicted instance, N_{GT} is the total number of ground truth instances, and $\mathbb{I}(\cdot)$ is the indicator function. This metric directly penalizes both missed teeth and spurious detections.

The tooth-level F1 score evaluates the model’s ability to correctly detect and segment individual teeth as distinct anatomical instances. It is computed based on connected-component matching between prediction and ground truth, with small components (volume < 11 voxels) discarded as noise. A ground truth tooth is considered correctly detected if at least 65% of its volume is covered by one or more predicted teeth, each of which contains no more than 70% non-overlapping (background) voxels. Conversely, a predicted tooth is deemed valid if it satisfies the symmetric matching condition with respect to the ground truth. The tooth-level F1 score is defined as the harmonic mean of tooth-level recall and precision:

$$F_1^{\text{tooth}} = \begin{cases} \frac{2 \cdot \text{Precision}_{\text{tooth}} \cdot \text{Recall}_{\text{tooth}}}{\text{Precision}_{\text{tooth}} + \text{Recall}_{\text{tooth}}}, & \text{if } \text{Precision}_{\text{tooth}} + \text{Recall}_{\text{tooth}} > 0, \\ 0, & \text{otherwise,} \end{cases} \quad (5)$$

where $\text{Recall}_{\text{tooth}} = \frac{\text{TP}_G}{|G|}$ is the ratio of correctly detected ground truth teeth (TP_G) to the total number of valid ground truth teeth ($|G|$), and $\text{Precision}_{\text{tooth}} = \frac{\text{TP}_P}{|P|}$ is the ratio of valid predicted teeth (TP_P) to the total number of valid predicted teeth ($|P|$). This metric emphasizes instance-wise detection accuracy rather than voxel-wise overlap, making it more clinically relevant for tasks such as individual tooth identification and counting.

In addition to the above metrics, The competition also considers the algorithm’s runtime and GPU memory consumption, ensuring that the runtime does not exceed 60 seconds in each case and that memory consumption is measured by the area under the memory–time curve. The evaluation process is thoroughly examined using these multi-dimensional metrics to ensure that the algorithms are reasonably evaluated in terms of both performance and efficiency.

The Challenge attracted research teams from several countries and regions around the world, with the largest number of participants from China. According to the participation statistics, Task 1 (2D tooth segmentation task) was participated in by 116 teams with a cumulative number of 346 submissions; Task 2 (3D CBCT tooth segmentation task) was participated in by 104 teams with a cumulative number of 158 submissions. These data reflect that researchers around the world have paid great attention to tooth image segmentation technology and extensively explored semi-supervised learning methods in real applications. To recognize outstanding contributions, cash prizes of \$500 were awarded to the top-performing team in each track, with additional souvenirs for other top teams presenting in person.

4. Methods

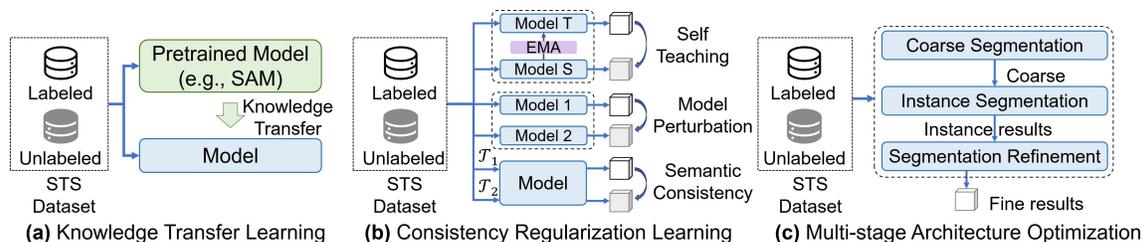


Figure 4: Overview of prominent methodological strategies employed by participants in the STS 2024 Challenge. The figure illustrates four key approaches: (a) Knowledge transfer with pretrained models, where pre-trained foundation models (e.g., SAM) are leveraged to improve segmentation. (b) Consistency regularization learning, including self-teaching, model perturbation, and semantic consistency (\mathcal{T}_1 and \mathcal{T}_2 denote two kinds of transformation). (c) Multi-stage architecture optimization decomposes the problem into multiple sub-problems and gradually obtained fine results.

Table 4: Summary of 2D submitted teams, model architecture, optimization strategy, and training settings.

Team	Model Architecture	Backbone	Optimizer	Loss Function	Device	Epochs	Batch Size
ChohoTech	YOLOv8	YOLOv8	Adam	CIoU, DFL, VFL	RTX 3090	100	32
Camerart2024	Self-Training Pipeline	DeepLabV3+	Adam	BCE, Dice	RTX 4090	200	4
Jichangkai	Two-Stage Semi-Supervised nnU-Net	nnU-Net	AdamW	Dice, CE	RTX 4090	150	4
Dew123	DICL Network	UNet	SGD	Dice, MSE, CE	RTX 3060	100	8
Junqiangmler	Semi-TeethSeg2024	VNet2d	AdamW	Dice, Cross-Entropy	RTX 4090	300	4
Isjinghao	SemiT-SAM	SAM	AdamW	Multi-component	RTX 3060	300	4
Lazyman	Cross Teaching Network	CNN + Transformer	SGD	Dice, MSE	RTX 4090	43	16
Caiyichen	YOLOv9	YOLOv9	Adam	CIoU, DFL, VFL	RTX 3060	100	16
Guo7777	ResUnet50 + SAM	ResNet50, SAM	Adam	BCEWithLogitsLoss, MSELoss	Tesla V100 -SXM2	300	4
Ccc2024	DAE-Net	Dual Attention Mechanism	Adam	Dice, IoU	RTX 4060 Ti	40	32

4.1. SAM-based Knowledge Transfer

In the dental segmentation task, major challenges include the scarcity of labeled data and the anatomical complexity of teeth. To address these issues, participants integrated the Segment Anything Model (SAM) to leverage its powerful general segmentation capabilities and edge sensitivity.

Chohotech combined YOLOv8 [8] and SAM [41] in a semi-supervised pipeline. YOLOv8 first generates bounding boxes to prompt SAM for fine pixel-level segmentation, creating high-quality pseudo-labels. These labels are refined via a quality filtering mechanism and used for iterative optimization. The method also adapts feature resolutions and loss functions for dental images and employs the Hungarian algorithm to resolve tooth number matching.

Guo777 proposed a framework integrating ResNet50 [46] and SAM-Med2D [42]. The ResNet50 encoder and U-Net decoder extract image features, while the SAM-Med2D module utilizes an attention mechanism to focus on key anatomical regions. To adapt to the intensity differences in medical X-rays, the input images were pre-

Table 5: Summary of 3D submitted teams, model architecture, optimization strategy, and training settings for STS MICCAI 2024 Challenge Task 2.

Team	Model Architecture	Backbone	Optimizer	Loss Function	Device	Epochs
Chohotech	3-Stage Pipeline	YOLOv8, U-Net	Adam	CIoU, DFL, VFL	NVIDIA A100	100
Houwentai	CFP 2-Stage Semi-sup. nnU-Net	nnU-Net	AdamW	Dice, CE	6 × RTX 4090	100
Madongdong	Semi-supervised YOLOv8	YOLOv8	Adam	CIoU, DFL, VFL	RTX 3090	300
Jichangkai	2-Stage nnU-Net Self-training	nnU-Net	AdamW	Dice, CE	RTX 3090	300
Junqiangmler	ROI Preproc. + VNet3d	VNet3d	AdamW	Dice, CE	RTX 4090	300
Gute_iici	2-Stage Unimatch	VNet (S1), Enc-Dec (S2)	AdamW	Unimatch	Tesla V100	100

processed using min-max normalization to rescale pixel values to $[0, 1]$ and resized to a consistent 512×512 resolution. This standardization ensures robust feature extraction despite varying lighting and exposure conditions. A semi-supervised strategy further expands the training set through pseudo-label generation and screening.

Isjinhao introduced SemiT-SAM, an innovative visual base model. Inspired by SAM-style architectures, it inherits MobileSAM’s lightweight ViT-Tiny backbone and produces multi-scale feature maps (C2–C5) through a simple feature pyramid. To handle the spatial variability of dental radiographs, images were resized while preserving the aspect ratio (max dimension 1024) and zero-padded to a fixed 1024×1024 input size, avoiding geometric distortion. In the semi-supervised phase, a teacher-student distillation strategy is used: the teacher generates pseudo-labels, which are filtered using class-confidence and mask-size thresholds to ensure high-quality supervision, following the SAM-based distillation workflow. The student model is jointly trained on labeled and pseudo-labeled data, optimizing the teacher weights via Exponential Moving Averaging (EMA) (Fig. 5).

4.2. Semantic Consistency Learning

Consistency learning enhances model robustness by ensuring feature and semantic alignment between labeled and unlabeled data, often through perturbations and multi-scale learning.

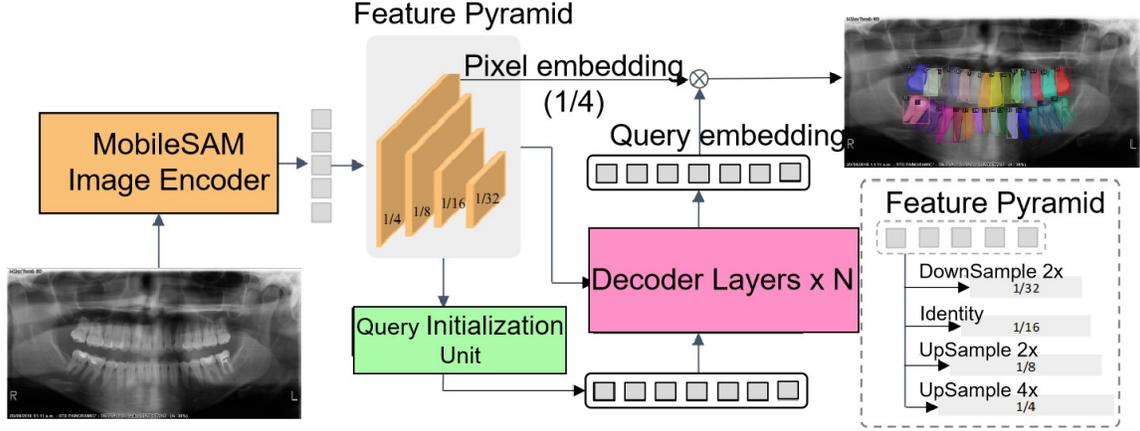


Figure 5: Architecture of the SemiT-SAM model, submitted by team 'Isjinhao' for the 2D challenge track. The model employs an encoder-decoder structure comprising an image encoder, a basic feature pyramid, a query initialization unit, and a mask decoder.

Houwentai proposed a Coarse-to-Fine Pseudo-labeling (CFP) method. An initial model trained on limited labeled data generates coarse pseudo-labels, from which high-confidence samples are filtered for fine segmentation training. The architecture improves upon the 3D nnU-Net [29] by incorporating residual modules for stability. Test-time augmentation (TTA) is applied during inference to integrate multiple predictions, significantly boosting accuracy and robustness under limited supervision.

Jichangkai developed a two-stage semi-supervised framework for both 2D and 3D tasks. In Stage 1, images are segmented into four anatomical quadrants using a low-resolution nnU-Net to simplify the spatial context. In Stage 2, each quadrant is processed by a full-resolution nnU-Net for detailed instance segmentation. A teacher-student framework iteratively generates and filters pseudo-labels, while specific mechanisms remove interfering data from adjacent quadrants. Finally, results are merged and renumbered to produce the complete segmentation (Fig. 6).

Camerart utilized a self-training approach for 2D segmentation. High-quality pseudo-labels are generated via multi-model integration and morphological manipulation, refined through 5-fold cross-validation. The DeepLabV3+ model employs a sigmoid function to handle tooth overlaps. During inference, segmentation is optimized using dual inference (original and flipped images) and contour extraction.

Lazyman combined U-Net [27] and Swin-UNet [47] in a cross-network co-training mechanism. U-Net captures local edge details, while Swin-UNet leverages Transformers for global context. The networks generate pseudo-labels for each other, enabling the model to learn from unlabeled data by enforcing consistency between local and

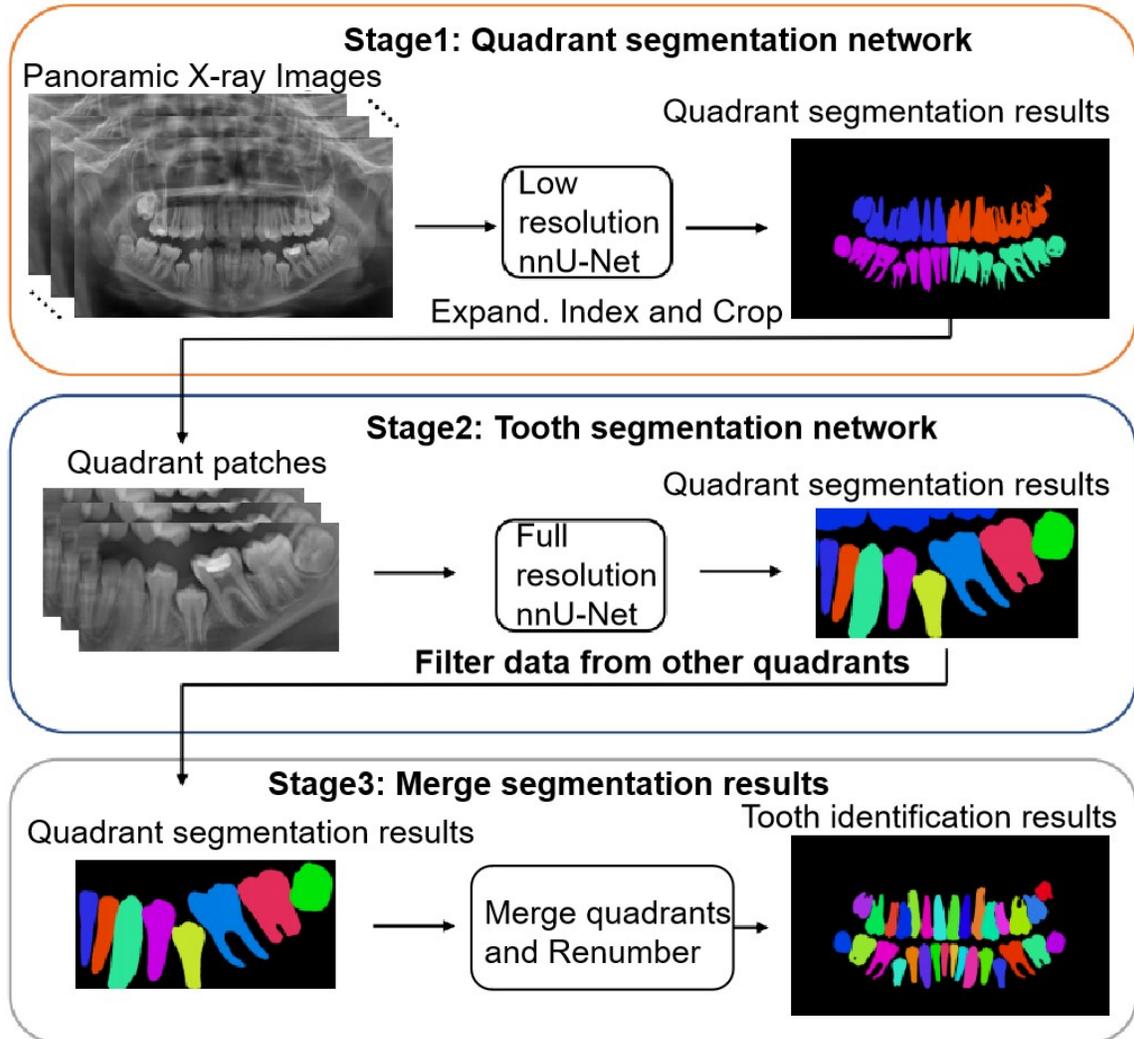


Figure 6: Schematic diagram of the two-stage segmentation method submitted by team 'Jichangkai' for the 2D challenge track. The first stage divides all teeth into four quadrants. The second stage identifies and segments each tooth in the quadrants. The third stage combines the results from all four quadrants to reconstruct the segmentation in the original image space.

global feature representations.

Dew123 proposed a deformable intrinsic consistency learning method. A semi-supervised stage uses a deformable convolutional module and cross-attention to generate pseudo-labels, followed by a fully-supervised stage combining these with labeled data. Semantic consistency is enforced via a composite loss function (Dice, cross-

entropy, and consistency loss).

Gute-iici applied the Unimatch framework in a two-stage 3D pipeline. A V-Net [28] first performs binary segmentation to extract ROIs. Subsequently, a multi-head network performs binary and multi-category segmentation. Feature perturbation at the bottleneck layer allows unenhanced data to supervise strongly enhanced outputs, improving feature learning from unlabeled data (Fig. 7).

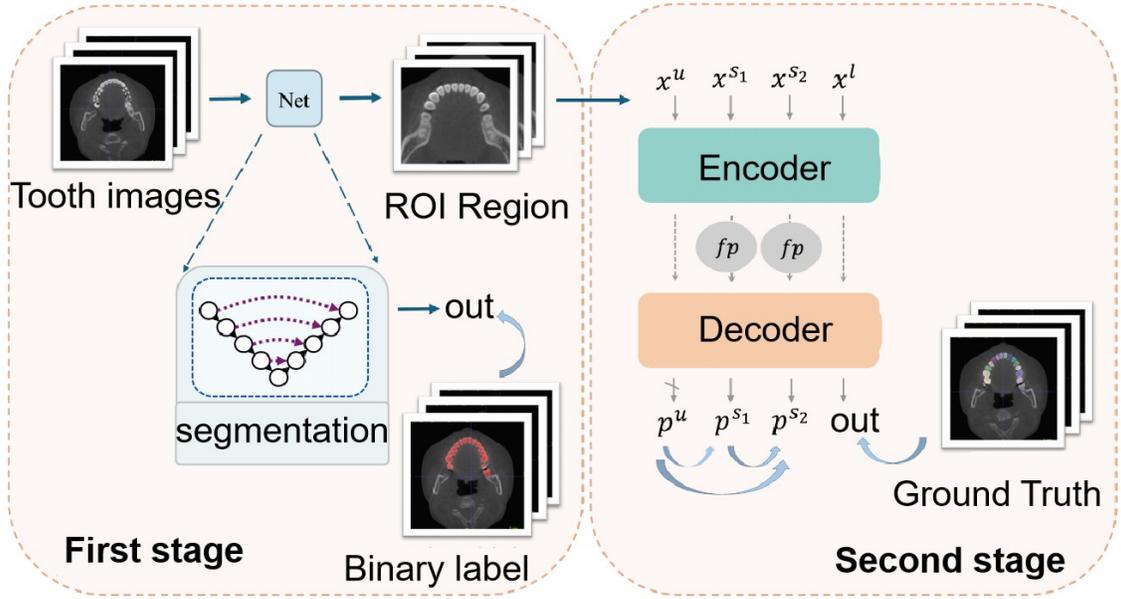


Figure 7: Two-stage semi-supervised tooth segmentation framework proposed by team 'Gute-iici' for the 3D challenge track. Unimatch is utilized for semi-supervised learning in the second stage.

4.3. Coarse-to-Fine Optimization

To handle the complexity of dental structures, several teams adopted multi-stage strategies that progressively refine segmentation from coarse localization to fine-grained detailing.

Haoyuuu proposed T3Net, a three-stage framework for CBCT images. Stage 1 uses a simplified Tiny V-Net [28] for coarse semantic segmentation to localize the ROI. Stage 2 employs a modified 3D ERFNet with spatial embedding, seed map, and prototype learning branches to generate coarse instance masks. Stage 3 refines these masks using a full-resolution Tiny V-Net. This cascaded approach effectively resolves under- and over-segmentation issues in complex 3D data (Fig. 8).

Madongdong optimized the YOLOv8 architecture [8] by introducing a coarse-to-fine structure, a decoupled head, and an anchor-free design. Their multi-stage

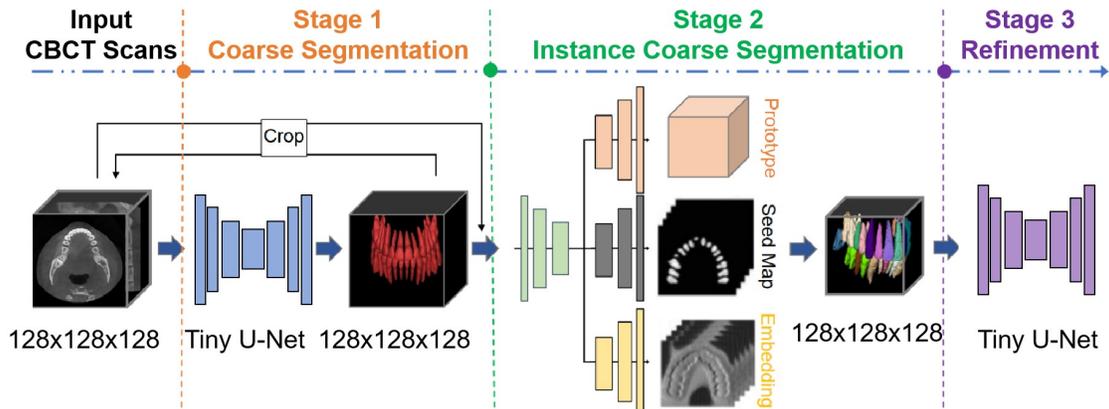


Figure 8: Overview of the T3 Net method, submitted by team 'Haoyuuuu' for the 3D challenge track. This automated pipeline performs instance-level segmentation and numbering of teeth and implants in CBCT images using a cascaded three-stage network.

training involves initially training on labeled data to generate pseudo-labels, which are then used to extend the dataset for a second stage of fine-tuning. This iterative process allows the model parameters to be optimized for high precision even with scarce labeled data.

Chohotech (3D) adapted their 2D strategy to CBCT data. The method begins with 2D Maximum Intensity Projections (MIP) processed by YOLOv8 for fast ROI detection. A modified 3D YOLOv8 then segments tooth instances within the ROI, using an adapted anchor strategy. Finally, a U-Net refines tooth boundaries using a combination of binary cross-entropy and Dice losses. This synergistic approach significantly improves efficiency and accuracy in volumetric segmentation.

4.4. Quantitative Evaluation

4.4.1. Quantitative Evaluation in 2D PXI Track

The comprehensive performance of the top ten teams in the 2D Panoramic X-ray (PXI) track, detailed in Table 6, underscores the advantage of semi-supervised learning (SSL) for this complex instance segmentation task. The most crucial finding of this challenge is the substantial performance gain achieved by SSL methods over a conventional supervised baseline within the constrained data regime of the challenge. We acknowledge that supervised learning remains the gold standard when large-scale annotations are available. However, in this specific scenario where the training set was restricted to just 30 labeled images (approx. 1% of the total dataset), SSL provided a critical performance boost. When trained only on the 30 labeled images, a fully-supervised nnU-Net baseline achieved an Instance Affinity (IA) of 44.17%. In

Table 6: Quantitative results for the 2D Panoramic X-ray (PXI) track. The top table shows the overall performance across multiple metrics including std in parentheses. The bottom table details the Instance DSC scores broken down by anatomical quadrant and patient age group (adult vs. child). The elevated standard deviations observed in certain teams may be attributed to the presence of a few cases with extreme complexity, as indicated in Fig. 9 and Fig. 16. **AUC_GPU** : Total GPU memory usage integrated over time, reflecting overall resource consumption. **F1**: The instance-level F1 score evaluates the balance between precision and recall for correctly detected individual tooth instances.

Team		Image-level		Instance-level				Algorithm-level	
Name	Ranking	DSC \uparrow	NSD \uparrow	DSC \uparrow	NSD \uparrow	IA \uparrow	F1 \uparrow	AUC_GPU (GB:s) \downarrow	RT (s) \downarrow
ChohoTech	1	92.03 (± 2.36)	95.89 (± 1.74)	83.59 (± 15.01)	87.52 (± 15.47)	88.53 (± 18.81)	92.54 (± 14.07)	7341.12 (± 1449.20)	13.29 (± 2.93)
Camerart2024	2	87.20 (± 4.58)	91.01 (± 4.70)	70.75 (± 16.87)	74.42 (± 17.56)	70.01 (± 21.71)	80.07 (± 18.44)	14250.98 (± 3870.70)	13.27 (± 3.67)
Jichangkai	3	92.01 (± 13.24)	94.87 (± 13.63)	78.98 (± 21.13)	81.93 (± 21.74)	80.06 (± 24.91)	85.98 (± 21.98)	25461.26 (± 1159.91)	55.90 (± 2.15)
Dew123	4	87.94 (± 4.69)	92.09 (± 4.53)	69.24 (± 21.58)	73.49 (± 22.42)	65.93 (± 25.71)	75.88 (± 23.94)	15088.50 (± 4448.06)	13.91 (± 4.49)
Junqiangmler	5	82.72 (± 10.13)	86.85 (± 10.53)	64.02 (± 18.43)	68.00 (± 19.27)	55.20 (± 24.15)	67.57 (± 23.56)	12483.38 (± 3398.61)	14.05 (± 3.96)
Isjinghao	6	82.64 (± 18.37)	86.34 (± 18.87)	64.71 (± 23.51)	67.75 (± 24.28)	67.31 (± 26.69)	76.68 (± 24.48)	27987.90 (± 3297.90)	21.13 (± 3.66)
Lazyman	7	59.76 (± 9.36)	87.05 (± 13.61)	49.22 (± 13.18)	72.60 (± 19.34)	8.57 (± 7.82)	14.50 (± 12.27)	13910.32 (± 5098.61)	11.81 (± 3.81)
Caiyichen	8	90.80 (± 2.53)	94.52 (± 2.40)	57.70 (± 27.60)	60.48 (± 28.82)	51.58 (± 33.84)	60.59 (± 34.32)	26666.57 (± 2129.66)	19.53 (± 1.16)
Guo77777	9	75.48 (± 9.44)	80.30 (± 9.51)	35.84 (± 13.04)	38.63 (± 13.61)	27.04 (± 16.43)	39.98 (± 20.82)	19694.26 (± 3798.62)	18.42 (± 3.18)
Cccc2024	10	91.59 (± 2.30)	95.15 (± 2.02)	26.60 (± 15.20)	27.75 (± 15.71)	15.38 (± 16.10)	24.09 (± 19.24)	17730.48 (± 1915.22)	13.46 (± 1.48)
baseline		89.38 (± 3.81)	51.61 (± 1.13)	71.12 (± 14.69)	41.45 (± 8.56)	44.25 (± 34.27)	55.54 (± 37.18)	11104.75 (± 4882.34)	1.22 (± 0.09)

stark contrast, the top-performing SSL method from team ChohoTech reached an IA of 88.53%, representing a relative improvement of over 100% (or an absolute gain of 44.36 percentage points). This result powerfully validates the central hypothesis of our challenge: that leveraging large amounts of unlabeled data through SSL is

Table 7: Quantitative results for the 2D Panoramic X-ray (PXI) track at quadrant level. The details include the Instance DSC scores broken down by anatomical quadrant and patient age group (adult vs. child).

Quadrant		Upper left		Upper right		Lower right		Lower left		Total
Team Name	Ranking	Adult	Child	Adult	Child	Adult	Child	Adult	Child	All age groups
ChohoTech	1	83.45 (±22.19)	84.87 (±19.58)	85.09 (±21.27)	81.77 (±25.44)	78.24 (±5.81)	76.25 (±12.23)	80.42 (±7.99)	80.35 (±6.34)	82.68 (±20.14)
Jichangkai	3	78.09 (±28.71)	79.09 (±26.14)	80.32 (±30.61)	78.32 (±30.64)	76.84 (±10.57)	74.72 (±13.27)	79.37 (±11.18)	77.30 (±11.44)	78.53 (±26.23)
Dew123	4	69.93 (±21.30)	70.77 (±25.55)	75.57 (±24.21)	73.27 (±26.93)	57.46 (±15.00)	57.06 (±16.47)	60.29 (±7.42)	61.67 (±10.04)	69.41 (±23.25)
Camerart2024	2	69.88 (±25.13)	71.06 (±22.95)	69.30 (±26.30)	68.63 (±24.94)	59.95 (±11.73)	69.78 (±13.80)	71.81 (±6.92)	66.34 (±13.28)	69.10 (±22.73)
Junqiangmler	5	62.07 (±21.95)	61.66 (±23.25)	67.55 (±24.86)	70.03 (±24.63)	61.67 (±17.35)	58.53 (±14.39)	68.53 (±9.34)	63.28 (±13.46)	64.81 (±22.20)
Isjinghao	6	64.59 (±26.93)	74.24 (±28.02)	66.79 (±31.65)	67.30 (±27.54)	46.23 (±20.96)	52.98 (±16.03)	49.95 (±16.62)	50.46 (±17.10)	64.12 (±27.84)
Caiyichen	8	60.33 (±30.55)	54.46 (±34.00)	47.40 (±35.58)	60.37 (±30.37)	49.72 (±23.86)	44.18 (±32.28)	41.36 (±28.38)	48.31 (±24.99)	53.45 (±32.27)
Lazyman	8	52.65 (±14.88)	52.20 (±17.58)	47.83 (±17.60)	50.64 (±17.98)	39.13 (±14.27)	40.40 (±10.01)	46.63 (±7.08)	44.04 (±8.44)	48.97 (±16.31)
Guo77777	9	48.59 (±24.99)	52.79 (±21.17)	43.49 (±22.69)	0.99 (±3.03)	43.04 (±17.74)	29.66 (±17.98)	27.04 (±12.91)	13.86 (±9.39)	34.66 (±27.00)
Cccc2024	10	38.63 (±21.17)	15.10 (±25.27)	6.49 (±8.79)	40.15 (±19.68)	37.39 (±14.32)	18.97 (±13.73)	18.62 (±17.19)	42.66 (±13.93)	26.06 (±23.36)
baseline		48.56 (±28.64)	36.37 (±23.09)	50.36 (±32.27)	38.99 (±26.25)	50.05 (±34.35)	48.08 (±27.04)	52.40 (±34.94)	48.35 (±25.35)	48.68 (±31.46)

essential for achieving high performance in instance-level dental segmentation.

The leading methods employed diverse yet effective SSL strategies. The top-ranked team, ChohoTech, utilized a YOLOv8-based detection-then-segmentation pipeline, achieving a robust instance-level Dice of 83.59%. Other leading approaches included a self-training framework (Camerart2024, rank 2) and a two-stage method built upon the powerful nnU-Net baseline (Jichangkai, rank 3). The raincloud plot in Fig. 9(a) visually confirms a clear performance stratification, with the top-tier teams consistently outperforming the lower-ranked participants in instance-level metrics, highlighting the efficacy of these advanced SSL frameworks.

The detailed quadrant-level analysis in Table 7 and the visualization in Fig. 9(e) show that most teams achieved relatively stable performance across the four anatomical quadrants (Q1: Upper Right, Q2: Upper Left, Q3: Lower Left, Q4: Lower Right).

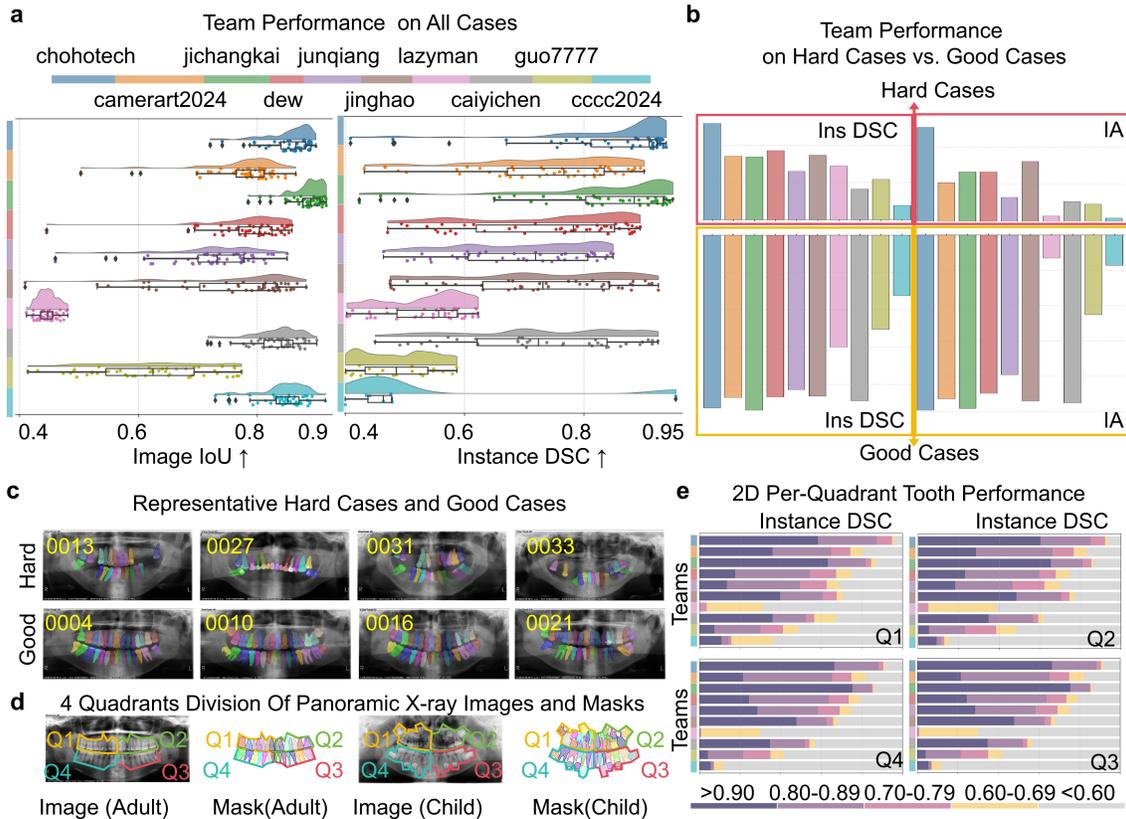


Figure 9: Analysis of 2D challenge track results. (a) Statistical summary of segmentation metrics for the top 10 teams on the test set. (b) Team performance on challenging cases and well-segmented cases. (c) Comparison of challenging cases with poor segmentation performance across most methods versus well-segmented cases, highlighting differences in tooth integrity. (d) Illustration of the four-quadrant division of a panoramic X-ray. (e) Instance DSC score distributions across four quadrants for participating teams. Quadrants: Q1 (Upper Right), Q2 (Upper Left), Q3 (Lower Left), Q4 (Lower Right). Ins DSC denotes Instance DSC.

This consistency is likely attributable to the inherent anatomical symmetry of human dentition, which aids model generalization. The analysis of challenging versus well-segmented cases in Fig. 9(b-c) further clarifies performance variations across "normal" and "abnormal" cohorts. As shown in Fig. 9 (b) and (c), well-segmented cases, such as 0004, 0010, 0016, and 0021, typically featured complete and well-aligned dentition. In contrast, challenging cases (e.g., 0013, 0027, 0031, 0033) were often characterized by severe tooth loss (edentulism) or complex pathologies. The performance on these "abnormal" cases degraded significantly for nearly all teams, as shown in Fig. 9 (b) by the lower and more variable scores.

Finally, the results highlight the nuanced trade-offs between different SSL strategies. For instance, the performance of the team Jichangkai reveals the importance of a multi-metric evaluation. Despite achieving high image-level and instance-level Dice scores, their comparatively lower IA score (80.06% vs. the winner’s 88.53%) suggests that while their nnU-Net based method produced spatially accurate masks, it struggled more with separating heavily overlapping teeth compared to the detection-first approach of ChohoTech. This underscores that different SSL architectures may excel at different aspects of the instance segmentation task—spatial accuracy versus topological correctness—providing valuable insights for future methods development.

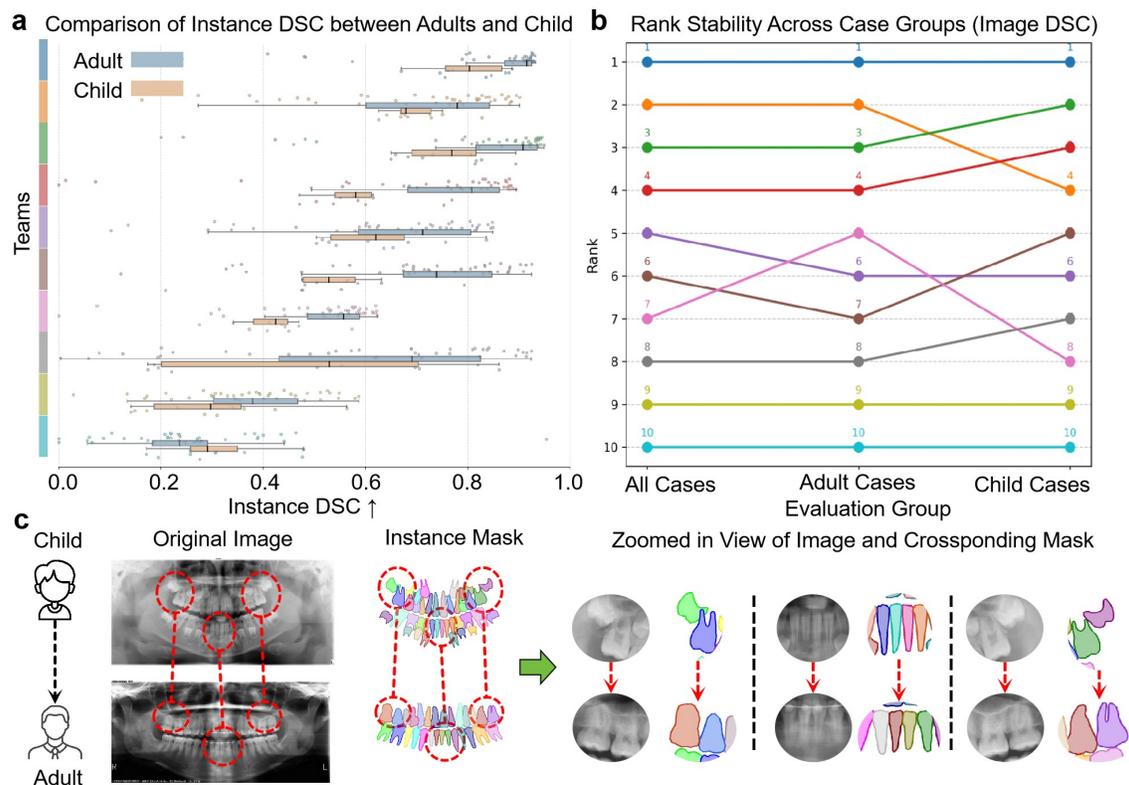


Figure 10: Performance analysis based on age groups in the 2D challenge track. (a) Statistical comparison of segmentation metrics for the top 10 teams on adult versus pediatric cases. (b) Illustration of anatomical differences and similarities between adult and pediatric dentition in Orthopantomogram. (c) Comparison of segmentation stability (e.g., variance in performance) across adult and pediatric cases for participating teams.

4.4.2. Performance on Challenging Pediatric and Adult Cohorts

To assess the robustness of the submitted SSL methods, we analyzed their performance when stratified by patient age (adult vs. child), a key challenge highlighted in this work. The results reveal that while pediatric cases are inherently more difficult, leading SSL methods demonstrated remarkable generalization capabilities.

As visualized in the raincloud plot in Fig. 10(a-b), a discernible performance gap exists between the two cohorts, with mean Instance DSC scores for adult cases generally being higher and less variable than for pediatric cases. This gap is attributable to both the smaller number of labeled pediatric examples in the training set and the intrinsic complexity of pediatric dentition. As shown in Fig. 10(c), pediatric anatomy often features smaller, more densely packed teeth with obvious overlap between deciduous and developing permanent teeth, making the instance segmentation task fundamentally more challenging. The rank stability analysis in Fig. 10(b) further underscores this challenge. To quantify the robustness of the rankings against demographic shifts, we calculated Spearman’s Rank Correlation Coefficient between team rankings on adult and pediatric cohorts. The analysis yielded a correlation coefficient of $\rho = 0.879$ ($p = 0.000814$) between adult and pediatric rankings, indicating strong overall consistency. While this demonstrates robustness of the top-performing methods across different age groups, several teams’ rankings did fluctuate significantly between the two cohorts, highlighting the challenge-specific performance variations. The correlations between overall rankings and age-specific rankings were $\rho = 0.964$ ($p = 7.32 \times 10^{-6}$) for the adult cohort and $\rho = 0.939$ ($p = 5.48 \times 10^{-5}$) for the pediatric cohort, confirming the stability of our evaluation framework. This indicates that their SSL strategies were not equally effective at generalizing across these distinct demographics. However, the quadrant-level results in Table 7 show that top-performing teams like ChohoTech and Jichangkai maintained impressive and relatively stable performance across all quadrants for both age groups. This indicates that their SSL strategies were not equally effective at generalizing across these distinct demographics. However, the quadrant-level results in Table 7 show that top-performing teams like ChohoTech and Jichangkai maintained impressive and relatively stable performance across all quadrants for both age groups.

Crucially, the success of the top methods on this difficult, underrepresented pediatric data is a significant outcome of this challenge. It validates the central premise that semi-supervised learning is a highly effective paradigm for developing robust models in clinical scenarios where acquiring labeled data for every subgroup is impractical. The ability to generalize from a few labeled examples to complex, unlabeled pediatric cases demonstrates a critical step towards building clinically viable and equitable AI tools.

Table 8: Quantitative results for the 3D Dental CBCT track. The top table shows the overall performance across multiple metrics. The bottom table details the Instance DSC and NSD scores for each anatomical quadrant. The increased standard deviations in certain teams may stem from the inclusion of a limited number of cases with exceptional complexity, as evidenced by the results in Fig. 16.

Team		Image-level		Instance-level				Algorithm-level	
Name	Ranking	DSC \uparrow	NSD \uparrow	DSC \uparrow	NSD \uparrow	IA \uparrow	F1 \uparrow	AUC_GPU (GB·s) \downarrow	RT (s) \downarrow
ChohoTech	1	93.54 (± 1.48)	97.42 (± 1.35)	92.15 (± 2.26)	96.60 (± 2.34)	98.39 (± 2.32)	99.56 (± 1.03)	233660.20 (± 119552.16)	60.76 (± 19.22)
Houwentai	2	76.79 (± 32.13)	78.52 (± 32.08)	83.59 (± 33.82)	73.77 (± 34.51)	75.24 (± 36.59)	99.38 (± 2.11)	829283.02 (± 308410.35)	210.37 (± 53.56)
Madongdong	3	83.57 (± 2.78)	80.87 (± 5.82)	77.87 (± 5.36)	73.63 (± 7.58)	88.18 (± 10.64)	93.63 (± 6.78)	48266.68 (± 19234.55)	52.82 (± 17.99)
Jichangkai	4	76.50 (± 29.64)	85.34 (± 19.77)	73.29 (± 32.22)	78.81 (± 28.20)	72.44 (± 39.12)	76.05 (± 37.21)	377331.12 (± 260326.42)	214.72 (± 110.78)
Junqiangmler	5	77.40 (± 26.22)	77.39 (± 27.61)	65.83 (± 31.04)	66.74 (± 31.67)	65.56 (± 36.29)	72.07 (± 36.98)	1004507.86 (± 528550.02)	114.11 (± 51.96)
baseline		71.99 (± 35.90)	73.34 (± 37.12)	30.80 (± 22.29)	37.28 (± 23.72)	14.98 (± 23.29)	18.04 (± 22.91)	5089073.15 (± 3453290.65)	130.21 (± 62.19)

Table 9: Quantitative results for the 3D Dental CBCT track at quadrant level. The details include the Instance DSC and NSD scores for each anatomical quadrant.

Quadrant		Upper left		Upper right		Lower right		Lower left		Total
Team Name	Ranking	Dice \uparrow	NSD \uparrow	Total Dice \uparrow						
ChohoTech	1	93.43 (± 1.79)	99.72 (± 0.75)	93.56 (± 1.76)	99.87 (± 0.44)	93.41 (± 1.72)	99.92 (± 0.14)	93.28 (± 1.70)	99.88 (± 0.52)	93.56 (± 1.49)
Houwentai	2	95.78 (± 2.84)	99.74 (± 1.59)	95.97 (± 1.13)	99.98 (± 0.05)	95.10 (± 2.77)	99.63 (± 1.12)	94.67 (± 4.11)	99.41 (± 2.34)	95.66 (± 1.47)
Madongdong	3	82.68 (± 3.12)	98.38 (± 2.55)	84.54 (± 2.75)	97.72 (± 2.11)	82.64 (± 5.62)	95.84 (± 5.99)	81.08 (± 4.31)	96.09 (± 3.37)	84.09 (± 2.92)
Jichangkai	4	77.08 (± 32.09)	85.57 (± 27.08)	71.43 (± 35.48)	81.78 (± 29.13)	72.67 (± 35.81)	78.11 (± 33.23)	75.12 (± 35.88)	81.17 (± 32.05)	76.54 (± 29.92)
Junqiangmler	5	79.94 (± 29.40)	86.00 (± 30.51)	80.12 (± 28.45)	86.22 (± 29.52)	69.49 (± 26.84)	76.74 (± 27.96)	73.55 (± 26.42)	81.71 (± 27.58)	77.42 (± 26.46)
baseline		36.17 (± 23.55)	40.76 (± 25.13)	32.05 (± 24.55)	38.51 (± 25.58)	34.85 (± 25.35)	43.18 (± 25.70)	33.60 (± 26.54)	39.04 (± 28.98)	71.98 (± 35.88)

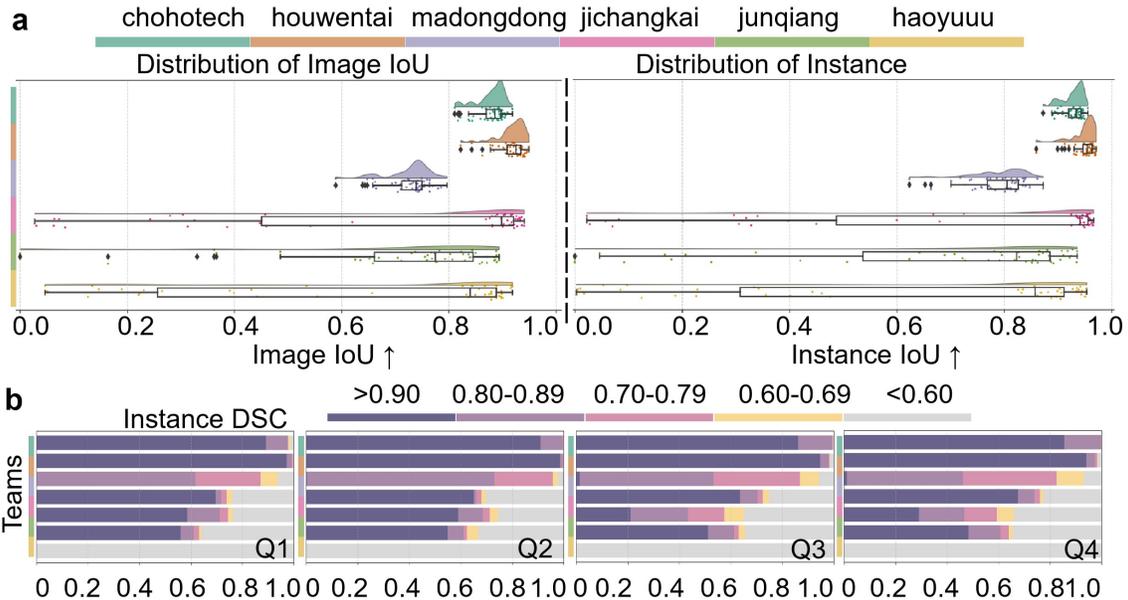


Figure 11: Summary and qualitative assessment of 3D challenge track results. (a) Statistical summary of test set metrics for the top 7 teams. (b) Distribution of Instance DSC scores for multiple teams across the four anatomical quadrants in CBCT.

4.4.3. Quantitative Evaluation in 3D CBCT Track

The 3D CBCT segmentation track posed a significant challenge, characterized by high-dimensional data and a high prevalence of imaging artifacts. This complexity was reflected in the final submissions, with only five teams successfully completing the task. The results, however, provide a powerful demonstration of semi-supervised learning’s efficacy in this demanding volumetric domain.

The most striking outcome of the 3D track is the dramatic performance improvement conferred by SSL. A fully-supervised 3D nnU-Net baseline, trained only on the 30 labeled CBCT scans, achieved a modest Instance DSC of 30.80%. In contrast, the winning SSL method from team ChohoTech achieved an Instance DSC of 92.15%. This represents an absolute gain of over 61 percentage points (a relative improvement of 197%), decisively validating the use of SSL for volumetric dental instance segmentation and proving its ability to overcome extreme data scarcity.

As detailed in Table 8 and visualized in Fig. 11(a), team ChohoTech’s approach, which adapted a YOLO-based pipeline to 3D, was not only the most accurate but also the most stable, evidenced by its low standard deviations across all metrics. This stability contrasts sharply with the performance of several other teams, including the second-place winner Houwentai, which exhibited very high variance (e.g., a

standard deviation of 33.82% for Instance DSC). This suggests that their SSL strategies, while effective on some cases, struggled to generalize across the full diversity of the unlabeled test set. This instability may indicate that their pseudo-labeling or consistency schemes were sensitive to domain shifts introduced by clinical factors like metal artifacts, causing degraded performance on unlabeled data that differed significantly from the small labeled set.

The quadrant-level analysis in Table 9 and Fig. 11(b) further reinforces the robustness of the top methods, which maintained consistently high performance across all four anatomical quadrants. The clinical relevance of these results is underscored by the excellent Normalized Surface Distance (NSD) scores achieved by the top teams. A high NSD score indicates the generation of segmentations with highly accurate surface boundaries, which is a critical prerequisite for clinical applications such as the digital design and fabrication of precise orthodontic appliances and surgical guides.

4.4.4. Comparative Analysis between 2D and 3D Tracks

A cross-task comparison reveals differences in performance, model complexity, and computational cost between the 2D PXI and 3D CBCT tracks, highlighting the distinct challenges of each modality. Four teams: ChohoTech, Jichangkai, Junqiangmler, and Madongdong competed in both tracks, allowing for a direct comparison of how their strategies adapted across dimensions.

The most striking difference lies in the potential for segmentation accuracy. As shown in Table 6 and Table 8, the winning team ChohoTech achieved a remarkable Instance DSC of 92.15% in the 3D track, substantially higher than their already impressive 83.59% in the 2D track. This suggests that the additional spatial information in 3D data is highly effective at resolving ambiguities inherent in 2D projections, most notably the superposition of tooth structures that complicates boundary delineation, leading to more accurate and robust instance segmentation once a model can effectively process volumetric data. However, capitalizing on this 3D advantage proved non-trivial. Other teams like Jichangkai and Junqiangmler saw a decrease in their average scores when moving from 2D to 3D, underscoring the challenge of successfully adapting SSL methods to a higher-dimensional space.

This difficulty is intrinsically linked to the algorithmic efficiency, where the computational burden of the 3D track was an order of magnitude greater. For ChohoTech, the average runtime (RT) increased from 13.29 seconds per 2D image to 60.76 seconds per 3D volume. More dramatically, their integrated GPU memory usage (AUC_GPU) surged from approximately 7,341 GB-s to 233,660 GB-s. This massive increase in resource requirements explains why fewer teams were able to complete

the 3D task and emphasizes a critical trade-off for clinical deployment: while 3D models can yield superior accuracy, their high computational cost may limit their practical application. The semi-supervised learning paradigm itself is also more demanding in 3D, as generating and refining pseudo-labels for volumetric data requires significantly more memory and computational power, a factor that contributes to the performance instability observed in several 3D submissions.

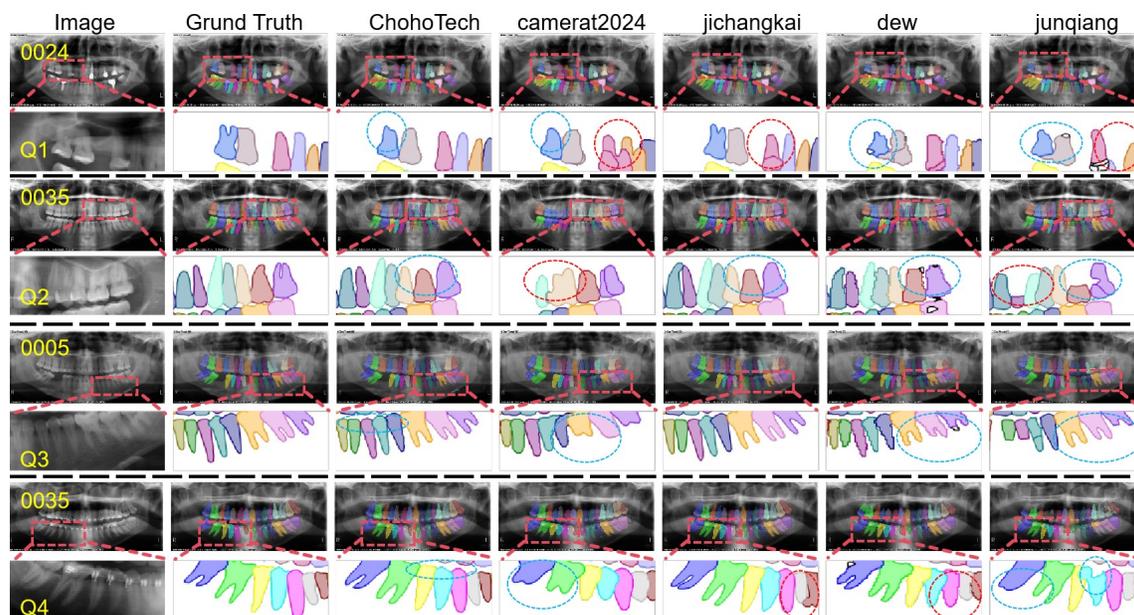


Figure 12: Qualitative comparison of 2D instance segmentation results from top-performing teams on representative adult panoramic X-ray (PXI) cases. Each case is presented with both a full view and a corresponding magnified quadrant view to facilitate detailed analysis of specific failure modes. These detailed views reveal common challenges, including the inaccurate segmentation of tooth apices (circled in blue) and the erroneous merging of overlapping teeth into a single instance (circled in red).

4.5. Qualitative Analysis of Segmentation Failures and Successes

To complement the quantitative metrics, a qualitative analysis of the segmentation results provides crucial insights into the specific clinical and anatomical challenges that current semi-supervised methods face.

4.5.1. Analysis of 2D PXI Segmentation

Visual inspection of adult cases in Fig. 12 reveals that while top-performing methods were proficient in segmenting well-defined teeth, two primary failure modes

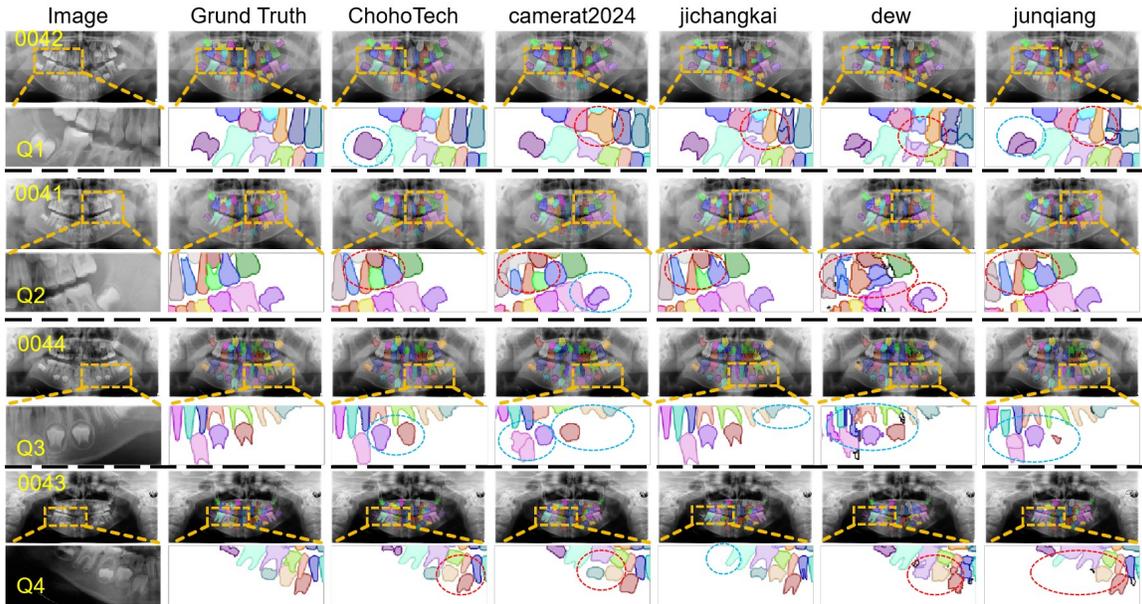


Figure 13: Qualitative comparison of 2D instance segmentation results from top-performing teams on representative pediatric panoramic X-ray (PXI) cases. Each case is presented with both a full view and a magnified quadrant view. Unlike adult dentition, pediatric cases present unique challenges due to the complex and irregular distribution of smaller teeth, leading to frequent and severe overlap. The magnified views highlight that most methods struggle in these dense, overlapping regions, often failing to accurately delineate individual small teeth, such as the erroneous merging of teeth (circled in red) or the inaccurate segmentation of developing apices (circled in blue).

emerged in challenging regions. First, the erroneous merging of overlapping teeth (circled in red) was a common issue. In 2D panoramic projections, the boundaries between crowded or superimposed teeth become ambiguous. This suggests that SSL models relying heavily on local pixel context can struggle to enforce instance separation without a strong prior, a problem that detection-first architectures are designed to mitigate. Second, many methods produced inaccurate segmentations of tooth apices (circled in blue). The apex is a small, low-contrast structure of high clinical importance for diagnosing periapical lesions. Its poor segmentation likely indicates that standard region-based loss functions (e.g., Dice) are dominated by the larger tooth crown, thus failing to penalize errors on these small but critical anatomical targets.

These challenges were significantly amplified in the pediatric cases shown in Fig. 13. Pediatric dentition, with its mix of smaller deciduous teeth and developing permanent teeth, creates a dense and highly irregular environment. The magnified

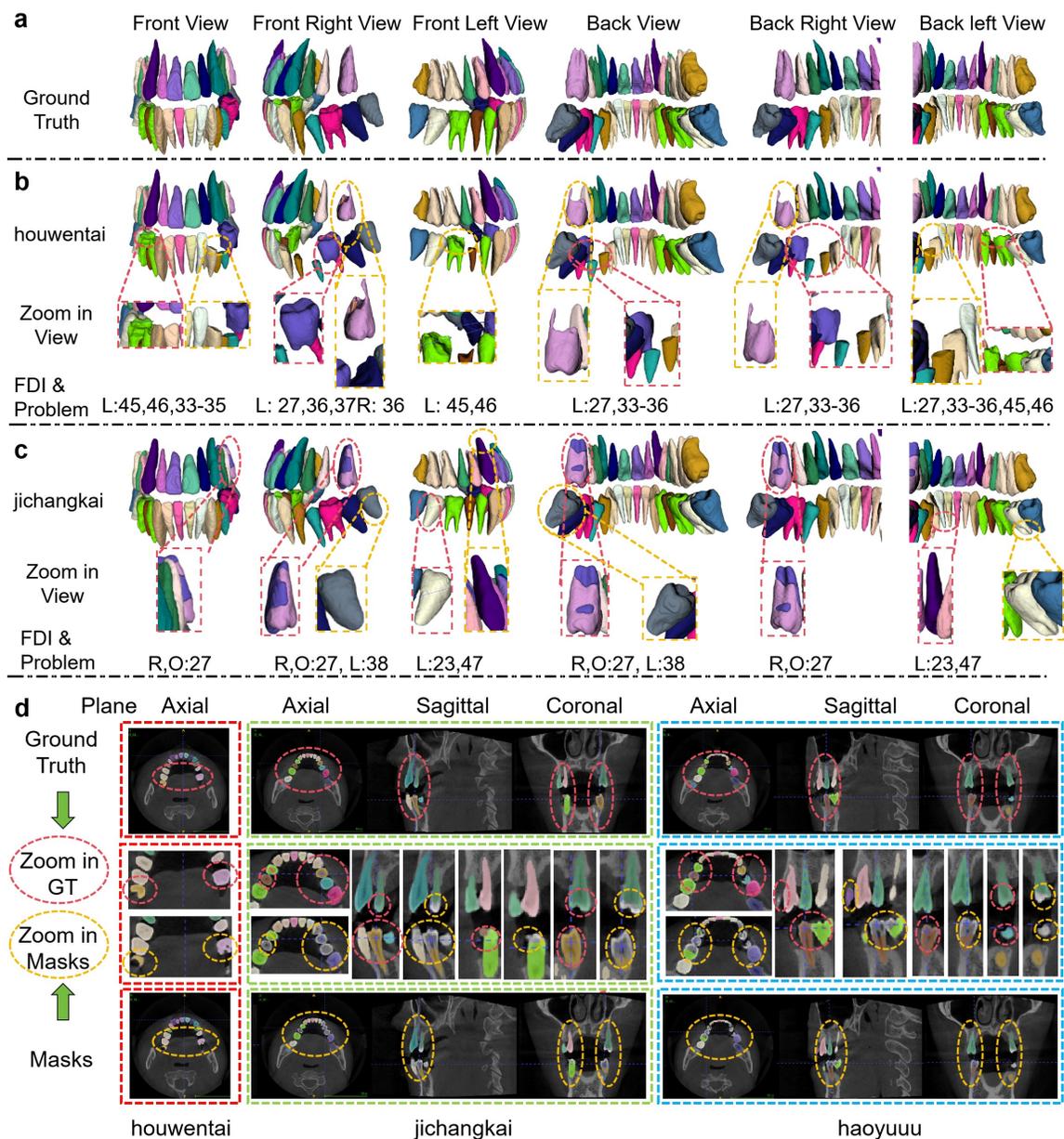


Figure 14: Visual comparison of 3D CBCT segmentation results. Multi-perspective views of the ground truth (a) are contrasted with outputs from top teams Houwentai and Jichangkai (b, c). Specific failure modes, including lack of segmentation (L), over-segmentation (O), and recognition errors (R) are indicated. A 2D slice comparison is provided in (d) for boundary detail.

views demonstrate that most methods struggled in these cluttered regions, often failing to delineate individual small teeth. This highlights a key difficulty for SSL: if the limited labeled data does not sufficiently represent the vast anatomical variability of pediatric development, the model may generate noisy or incomplete pseudo-labels for unlabeled pediatric cases, hindering its ability to learn robust features for these challenging objects.

4.5.2. Qualitative Analysis of 3D Tracks at zoomed in view

For the 3D track, Fig. 14 illustrates typical error modes in volumetric segmentation. While leading methods successfully reconstructed the 3D morphology of most teeth, three primary failure patterns emerged: under-segmentation (e.g., merging adjacent teeth, labeled 'L'), over-segmentation (e.g., erroneous extensions into surrounding bone, labeled 'O'), and recognition errors (e.g., correct segmentation but incorrect tooth identification, labeled 'R'). These qualitative observations corroborate the quantitative findings, underscoring that accurate instance separation and identification remain critical challenges in 3D CBCT segmentation.

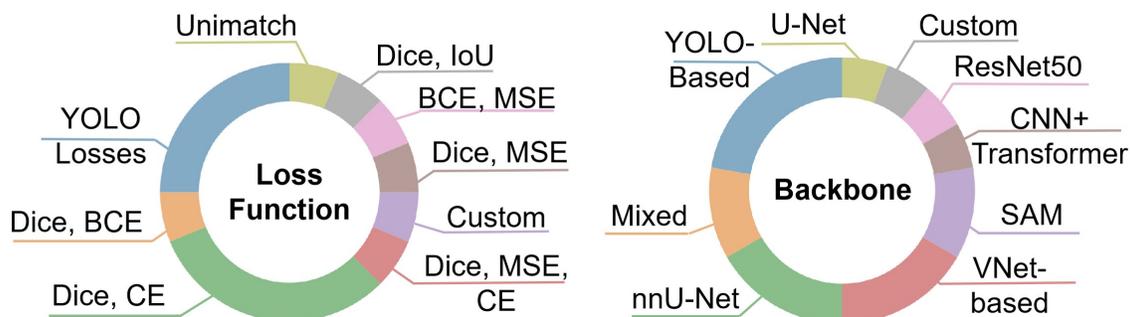


Figure 15: Distribution of loss functions employed and backbone architectures by participants in the 2D and 3D challenge tracks. (Left) Analysis of loss function usage indicates a strong preference for composite loss functions over single-objective training. The combination of a region-based loss (Dice) and a pixel-wise loss (Cross-Entropy) was the most common strategy. (Right) Breakdown of backbone model choices, revealing that nnU-Net and V-Net were the most prevalent architectures. A notable trend was the integration of the Segment Anything Model (SAM) as a foundational component.

4.6. Methodological Insights from Leading Teams

An analysis of the methods submitted by the top-performing teams in both the 2D and 3D tracks reveals a clear consensus on several effective strategies, from architectural choices to training paradigms. While a diversity of approaches was observed,

the most successful solutions consistently converged on a few key principles that effectively addressed the challenges of semi-supervised instance-level tooth segmentation. The distribution is depicted in Fig. 15.

4.6.1. Network Architectures

A dominant trend was the adoption of robust, highly specialized backbone architectures. The nnU-Net framework, renowned for its automated pipeline configuration and strong performance in medical imaging segmentation, emerged as a cornerstone for several leading teams, including Jichangkai and Houwentai. Its prevalence underscores the value of a well-tuned, powerful baseline in competitive challenges. Concurrently, some teams, such as ChohoTech and Caiyichen, successfully employed detection-based models like YOLOv8 and YOLOv9 as a foundational step. This indicates a strategic preference for two-stage "detect-then-segment" pipelines, where initial tooth localization simplifies the subsequent fine-grained segmentation task. Classic architectures like V-Net and U-Net also remained relevant, often adapted for specific tasks.

4.6.2. Semi-supervised Techniques

Beyond the choice of backbone, a key determinant of success was the implementation of multi-stage, coarse-to-fine refinement strategies. Rather than relying on a single end-to-end model, top teams often decomposed the complex problem into more manageable sub-tasks. The YOLO-based approaches are a prime example, using object detection for coarse localization before a dedicated segmentation head provides the final mask. Similarly, the winning method from Jichangkai implemented a two-stage process that first divided the image into quadrants and then performed high-resolution segmentation within each, effectively managing complexity and improving precision. This pattern suggests that explicitly guiding the model from a high-level overview to fine-grained detail is a highly effective strategy for dental instance segmentation.

4.6.3. Loss Function

Regarding the training objectives, no single loss function was universally dominant; instead, the clear trend was the use of composite loss functions. The most common combination for segmentation-focused models was a weighted sum of a region-based loss, typically the Dice loss, and a pixel-wise classification loss, such as Cross-Entropy (CE) or Binary Cross-Entropy (BCE). This hybrid approach balances the need for accurate spatial overlap (Dice) with correct pixel-level classification (CE/BCE), a standard practice that proved effective here. Teams employing

YOLO-based backbones naturally utilized their corresponding specialized loss functions, such as CIoU, DFL, and VFL, tailored for detection and localization accuracy.

4.6.4. Large-scale Foundational Models

Finally, a notable emerging trend was the integration of large-scale foundational models, particularly the Segment Anything Model (SAM) and its medical-specific variant, SAM-Med2D. Teams like Isjinghao and Guo7777 leveraged these models, likely to generate high-quality pseudo-labels, provide strong feature initialization, or act as a powerful component within their segmentation pipeline. This highlights a strategic shift towards harnessing the powerful prior knowledge embedded in foundation models to enhance performance in data-scarce, semi-supervised scenarios. In summary, the leading solutions in the STS 2024 challenge skillfully combined robust backbones with multi-stage refinement pipelines, optimized them with composite loss functions, and, in several cases, augmented their performance by integrating knowledge from foundational models.

4.6.5. Analysis of Limitations and Failed Approaches

In contrast to the success of multi-stage and foundation-model-based approaches, an analysis of lower-performing entries reveals that the lack of explicit anatomical priors was a primary cause of failure. Methods relying solely on pixel-level consistency regularization often failed to capture the global topological arrangement of the dentition. Specifically, in 2D OPG tasks, the absence of 'horizontal' spatial constraints (i.e., the relative ordering and separation of teeth along the dental arch) frequently led to the merging of adjacent, overlapping teeth. This failure mode was effectively mitigated in top-tier methods by using detection-based priors (e.g., YOLO bounding boxes) that enforce instance separation before segmentation. Similarly, in the 3D CBCT track, approaches that processed data without sufficient 'vertical' inter-slice consistency struggled to maintain volumetric integrity. Models lacking 3D context often produced fragmented segmentations, particularly for fine structures like tooth apices, as they failed to leverage the continuity of tooth anatomy across the Z-axis. Furthermore, standard self-training pipelines without robust noise-filtering mechanisms (such as uncertainty estimation or ensemble voting) tended to overfit noisy pseudo-labels. This was particularly evident in 'hard' cases involving tooth loss or implants, where the domain shift in the unlabeled data led to error propagation, limiting the models' generalization capability."

4.7. Performance Across Anatomical Regions and Demographics

To further dissect performance, we analyzed Instance DSC scores in different anatomical quadrants and demographic groups, as shown in Tables 6 and 8. For

the 2D PXI track, the results were stratified by quadrant and by patient age (adult vs. child). This detailed breakdown confirms the observations from the qualitative analysis: on average, segmentation performance on pediatric cases was slightly lower than on adult cases across most teams and quadrants. For example, the top team, Chohotech achieved a total Instance DSC of over 82%, with relatively stable performance across both adult and pediatric groups, demonstrating the robustness of their semi-supervised approach. This result is significant, as it confirms that SSL can be effectively applied to challenging pediatric data where labeled examples are particularly scarce, thereby addressing a key goal of this challenge. In the 3D CBCT track, performance across quadrants was generally high for the top teams, with Houwentai and Chohotech achieving Instance DSC scores above 93% in all regions. The lower-ranked teams showed more variability, indicating less consistent performance across the full dental arch.

4.8. Algorithm Efficiency and Clinical Relevance

The clinical applicability of AI models depends not only on accuracy but also on computational efficiency. Table 6 and 8 provides metrics for GPU memory consumption (AUC_GPU) and runtime (RT). In the 2D track, the top-ranked method from ChohoTech was not only accurate but also efficient, with a runtime of approximately 13 seconds per image. Several other methods also achieved inference times under 20 seconds, suggesting a strong potential for integration into clinical workflows where a near-real-time response is desirable. In contrast, team Jichangkai’s method, while highly accurate, required a significantly longer runtime (55.90s), which might be less suitable for interactive applications. The computational demands for the 3D track were substantially higher. The fastest method (Madongdong at 52.82s) was still considerably slower than the 2D methods. Other top-performing 3D models required several minutes for inference, and their GPU memory consumption was an order of magnitude higher than their 2D counterparts. This starkly illustrates the trade-off between the detailed anatomical information provided by 3D CBCT and the computational resources required to process it, a critical consideration for deployment in real-world clinical settings.

5. Discussion

5.1. Main findings

Dental image segmentation presents unique challenges that differentiate it from other medical imaging domains. The STS 2024 Challenge underscored several critical

difficulties inherent to dental datasets: 1) Limited Segmentation Accuracy of Previous Methods; 2) Inadequate Integration of Segmentation and Tooth Position Recognition; 3) High Proportion of Unannotated Data. Traditional methods rely heavily on fully supervised learning paradigms, which demand extensively annotated datasets. However, annotating dental images is labor-intensive and requires specialized expertise, leading to a scarcity of high-quality labeled data. Consequently, previous approaches exhibited limited precision in accurately delineating tooth boundaries, especially in cases involving overlapping structures or varying image qualities. The complexity of dental anatomy, with its intricate and closely packed tooth structures, further exacerbates segmentation inaccuracies.

Beyond mere segmentation, accurate tooth position identification is crucial for comprehensive dental analysis and clinical decision-making. Prior methodologies typically treated segmentation and tooth position recognition as separate tasks, leading to suboptimal performance in integrated scenarios. This discrimination fails to leverage the inter-dependencies between segmentation and positional data, resulting in inconsistencies and reduced overall accuracy. The inability to concurrently address segmentation and tooth position recognition limits the utility of automated systems in practical dental applications.

The STS 2024 Challenge introduced several innovative elements in dataset construction and challenge design to address the aforementioned challenges. First, to mitigate the scarcity of annotated data, we meticulously curated a dataset encompassing 2D panoramic X-ray images and 3D CBCT tooth volumes. This multimodal dataset enriches the dataset and facilitates the development of algorithms capable of handling different imaging modalities, encouraging the creation of more robust and generalizable models.

Furthermore, including pediatric and adult dental images ensures that models can adapt to variations across age groups, enhancing their applicability in real-world clinical settings. Our challenge incorporated a multi-faceted evaluation framework that assesses participants' algorithms on several levels, ensuring a comprehensive performance evaluation: Instance-Level Evaluation, Image-Level Evaluation.

The results from the STS 2024 Challenge reveal several noteworthy insights: Top-performing teams predominantly leveraged semi-supervised learning frameworks, effectively utilizing the limited labeled data alongside the abundant unlabeled data. Techniques such as integrating pre-trained models like the Segment Anything Model (SAM), consistency regularization learning, and multi-stage architecture optimization significantly enhanced segmentation accuracy. This trend aligns with the broader trend in medical imaging, where semi-supervised methods are increasingly recognized for their ability to overcome data scarcity.

In the 3D CBCT segmentation task, multi-stage training strategies yielded superior performance. Participants could incrementally improve segmentation precision and robustness by sequentially refining the model through multiple training phases. This approach facilitates better feature learning and mitigates the risk of overfitting, particularly in complex 3D structures.

Despite the overall success of semi-supervised methods, a notable portion of participants did not employ any semi-supervised strategies. This variability reveals a critical gap between the demonstrated potential of SSL and its practical adoption, which have demonstrated clear advantages in handling unannotated data.

Across the diverse methodologies, participants predominantly introduced perturbations at multiple levels to augment model training and generalization: **Data-Level Perturbations:** Techniques such as weak-strong augmentations and the addition of Gaussian noise were employed to create varied input data representations, fostering robustness against input variations. **Model-Level Perturbations:** The integration of heterogeneous network architectures (e.g., CNNs and Transformers) facilitated implicit consistency regularization, promoting diverse feature learning and reducing model-specific biases. **Training Cycle Perturbations:** Frameworks like the teacher-student model incorporated temporal perturbations by iteratively refining pseudo-labels, thereby enhancing feature learning and mitigating overfitting.

In terms of consistency constraints, methods varied across three primary dimensions: **Soft Constraints:** Utilization of soft logits and features as supervisory signals encouraged smooth decision boundaries and feature consistency across different perturbations. **Hard Constraints:** Binary segmentation outputs were enforced through hard label assignments, ensuring definitive delineation of anatomical structures. **Structural Constraints:** Incorporation of edge detection operators (e.g., Sobel) imposed structural integrity on anatomical boundaries, aligning segmentation outputs with inherent anatomical features.

Building upon the insights gained from the STS 2024 Challenge, several avenues for future research and development emerge. First, efforts should be directed toward creating more extensive and diverse annotated dental datasets encompassing various anatomical variations, imaging modalities, and clinical conditions. Collaborative initiatives and data-sharing platforms can facilitate the accumulation of comprehensive datasets necessary for training robust deep-learning models.

Further exploration of semi-supervised and unsupervised learning methods can enhance model performance and reduce dependency on labeled data. Techniques such as self-supervised learning, transfer learning, and active learning hold promise for improving segmentation accuracy and efficiency. Developing models that concurrently address segmentation and tooth position identification can yield more holistic

and clinically relevant outputs. Multi-task learning frameworks can enhance the synergy between tasks, improving overall performance and utility. Bridging the gap between research and clinical practice requires rigorous validation of segmentation models in real-world settings. Future studies should focus on integrating automated segmentation tools into clinical workflows and assessing their impact on diagnostic accuracy, treatment planning, and patient outcomes.

Enhancing the computational efficiency of semi-supervised models is crucial for their adoption in clinical environments. Research should optimize models for real-time processing without compromising segmentation accuracy, ensuring their practicality and scalability. Subsequent iterations of the STS Challenge can incorporate more diverse and complex datasets, introduce additional evaluation metrics that reflect clinical utility, and encourage the development of lightweight models suitable for deployment in various clinical settings. Such enhancements can further drive innovation and accelerate the adoption of advanced segmentation techniques in dentistry.

5.2. Case Difficulty and Metric Correlation

The case difficulty and metric correlation analysis in Fig. 16 offers valuable insights into the dataset’s intrinsic challenges and the robustness of our evaluation framework.

The difficulty analysis in Fig. 16(a, c) reveals that the primary challenges differed by modality. In the 2D PXI track, the most difficult cases (e.g., Case 0031, which had the lowest scores) were those with severe pathologies like extensive tooth loss. This is a classic problem for semi-supervised learning, as these atypical cases deviate from the primary distribution of the unlabeled data, hindering the generation of effective pseudo-labels. This suggests a need for SSL methods that are more robust to domain shift within the dataset itself. In contrast, 3D CBCT challenges were often related to acquisition variations and fine anatomical details. For instance, Case 0024 was identified as the most challenging 3D sample. Visual inspection revealed that this volume suffered from a severe scan position offset, where the dentition was shifted entirely to the anterior third of the image space. This spatial anomaly caused failure in models relying on standard centering priors. This points to architectural limitations in current models, which may require more specialized designs to capture fine-grained features in volumetric data.

The strong positive correlation between different evaluation metrics, as shown in the matrices from Fig. 16(b, d), is an interesting finding. The high concordance across metrics like Instance DSC and Instance Affinity (IA) indicates that the team rankings were stable and not dependent on a specific metric choice. This validates the robustness of our evaluation framework and suggests that the top-performing

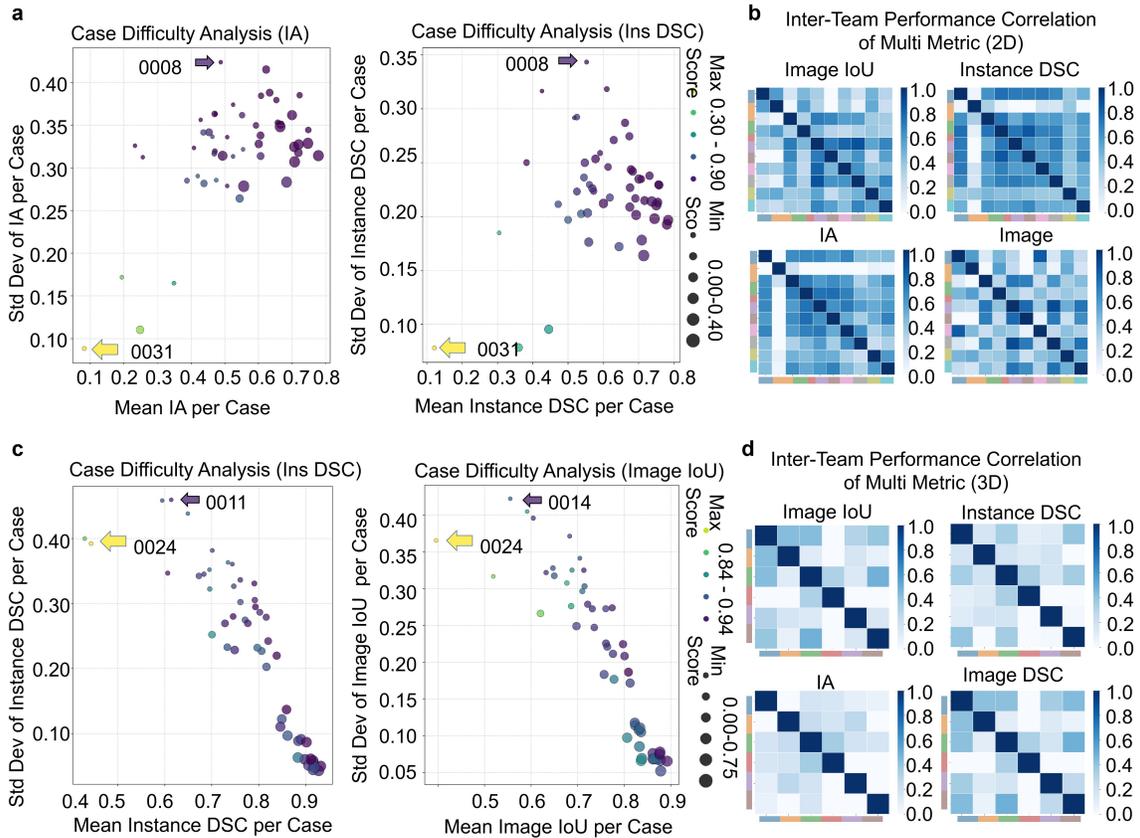


Figure 16: Case difficulty and inter-metric correlation analysis for the 2D (a, b) and 3D (c, d) tracks. (a, c) Case Difficulty Analysis: Plots identifying the most challenging cases (0031 in 2D, 0024 in 3D). Difficulty in 2D was linked to severe tooth loss, while in 3D, it arose from small targets (e.g., apices) and inter-tooth misidentification. (b, d) Inter-Metric Correlation: Matrices for both tracks showing strong concordance across metrics, indicating consistent team rankings.

methods were genuinely superior across multiple criteria of segmentation quality.

5.3. Human-Machine Collaborative Annotation Study

To quantitatively evaluate the clinical utility of semi-supervised segmentation algorithms in this challenge, we conducted a human-machine collaborative annotation study. The study involved two dentists (a junior clinician with more than two years of experience and a senior specialist with more than ten years of experience). The dentists annotated a curated test set of 30 dental scans, including 10 X-rays from child, 10 X-rays images from adult, and 10 CBCT volumes from adult. To simulate a realistic clinical workflow, dentists utilized LabelMe for 2D X-ray polygon anno-

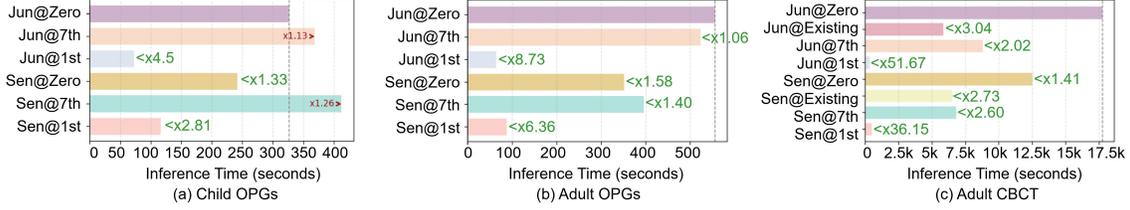


Figure 17: Bar plots of X-ray and CBCT labeling time under varied methods. The **green arrow** indicates a time reduction (acceleration) and the **red arrow** an increase (deceleration) compared to the baseline (typically the Junior@Zero). Jun: Junior dentist; Sen: Senior dentist.

tation and 3D Slicer for 3D CBCT volumetric refinement. These industry-standard tools were used across all strategies to ensure consistent time measurement. The study follows four distinct strategies: (i) Full manual annotation (i.e., *From scratch*); (ii) Correction of predictions from existing methods (only available for the CBCT); (iii) Correction of 7th-ranked outputs (relatively low-quality); and (iv) Correction of 1st-ranked outputs (relatively high-quality).

Annotation time per case was rigorously recorded to assess workflow efficiency, as shown in Fig. 17. This study focuses on three key dimensions: (i) the relationship between algorithmic performance (as reflected by model rankings) and time saved during clinical workflows; (ii) modality-specific challenges in 2D vs. 3D tasks, including anatomical variability across pediatric and adult populations; and (iii) how clinician expertise interacts with automated tools to optimize workflow efficiency.

5.3.1. Dramatic Time Savings in Clinical Workflows

The top-performing semi-supervised model (the 1st-ranked team) demonstrated transformative efficiency gains across all modalities. In 2D X-rays (a historically unaddressed domain), the labeling efficiency improved 4.5 \times for the junior dentists (i.e., from 326.10s to 72.50s per case) and 2.11 \times for the senior (i.e., from 244.50s to 116.10s). Critically, the 1st-ranked model outperformed existing CBCT methods by 94.1% (i.e., the junior’s time that decreased from 5841.80s to 343.30s), validating its superiority in handling sparse-annotation scenarios. For the CBCT volumes, the annotation time plummeted from 3.5 4.9 hours to 5.7 8.2 minutes per volume, which is clinically critical for efficient orthodontics and implant planning.

Moreover, the segmentation quality directly impacts the annotation efficiency. The relatively low-quality initial predictions may increase the cognitive load and correction time of dentists. For instance, the 7th-ranked model prolonged senior dentists’ child X-ray annotation by 68.5% (i.e., increasing from 244.50s to 412.00s), indicating that subpar segmentation disrupts clinical efficiency. Conversely, the 1st-

ranked model reduced time while improving consistency, where the standard deviations decreased across tasks. This highlights that only high-precision algorithms enable reliable human-machine collaboration, which is one of the core goal of our semi-supervised challenge.

Notably, the CBCT modality that is essential for complex orthodontic procedures, showed the greatest absolute time savings (e.g., from 4.9 hours to ≤ 0.1 hour / volume). This aligns with our challenge’s focus on underrepresented 3D data, addressing a critical barrier in dental area. For adult X-rays, the 1st-ranked model reduced annotation time by 88.2% (e.g., the junior’s time decreased from 555.50s to 63.60s), outperforming the 7th-place model by 87.7%. Notably, for pediatric X-rays—a domain previously lacking public datasets—the methods achieved near-adult levels of efficiency (e.g., 72.5s vs. 65.7s for the junior), proving semi-supervised methods can overcome pediatric data scarcity.

5.3.2. Implications for Clinical Deployment

The junior dentist benefited disproportionately from high-quality AI assistance. In pediatric X-rays, their efficiency gain with the 1st-ranked model ($4.5\times$) exceeded Seniors’ ($2.1\times$), narrowing the experience gap: the junior-senior time differential dropped from 81.6s (manual) to 43.6s (AI-assisted). However, Seniors exhibited lower tolerance for imperfect predictions—when using the 7th-place model, their pediatric X-ray annotation time surged by 68.5%, versus only 12.9% for Juniors. This implies that AI assistance most empowers early-career clinicians, while experts require near-perfect segmentation for workflow integration. Based on the measured throughput (decreasing from 4 hours to 5 minutes per volume), we estimate that a Senior dentist could process approximately 96 CBCT volumes/8-hour workday with 1st-ranked model assistance versus 2 volumes/workday manually—enabling large-scale data curation for AI development. Combining Junior clinicians with state-of-the-art AI achieved 73% of Senior-level efficiency (e.g., pediatric X-rays: 72.5s vs. Senior’s 116.1s) at lower resource cost. This validates our challenge’s design premise: semi-supervised learning bridges data scarcity in specialized domains (pediatrics, CBCT) while optimizing clinical resource utilization.

In a nutshell, our challenge’s top semi-supervised model reduced dental annotation effort by 52–98%, demonstrating that human-AI collaboration—when powered by high-precision initial segmentation—can overcome historic barriers in pediatric and 3D dental imaging, accelerating the translation of AI tools into clinical practice. Our challenge catalyzed algorithms that transform dental image annotation from a hours-long manual task to a minutes-long collaborative effort, with profound implications for clinical scalability and global dental health equity. Future work should focus

on real-time interactive refinement tools to further harness human-AI synergies. The results underscore the critical role of semi-supervised algorithms in enabling scalable, cost-effective dental diagnostics while maintaining clinical accuracy—a prerequisite for widespread adoption in resource-constrained settings. Also, this study directly addresses two critical gaps in dental AI research: (1) the absence of benchmarks for human-AI collaboration efficiency in pediatric dental imaging, and (2) the lack of comparative data on how semi-supervised model quality impacts real-world clinical labor. By simulating real-world refinement scenarios, we quantify how algorithmic performance translates to tangible clinical time savings and evaluate the viability of semi-supervised solutions for underserved tasks (e.g., pediatric X-rays and CBCT segmentation, where public datasets are absent).

5.4. Limitations and Future Directions

While the STS 2024 Challenge made strides in advancing dental image segmentation, it has limitations. The dataset, although diverse, remains limited in size and scope, potentially affecting the generalizability of the findings. Additionally, the challenge primarily focused on segmentation accuracy, with less emphasis on other critical factors such as processing speed and integration with clinical workflows. Furthermore, while representative of real-world scenarios, the high proportion of unannotated data may have introduced biases that could influence the performance of semi-supervised algorithms.

Looking ahead, the trajectory of this field must shift from algorithmic competition to clinical translation. A critical open question is whether current performance levels are sufficient to support high-precision autonomous systems, such as AI-guided robotic dental implant surgery, where error margins are measured in fractions of a millimeter. To answer this, future research should move beyond geometric metrics (like Dice) and incorporate downstream functional benchmarks—evaluating, for instance, how segmentation errors impact the stability of orthodontic movement simulations or the fit of 3D-printed surgical guides. Furthermore, investigating the theoretical 'upper bound' of performance with fully labeled, multi-expert consensus datasets is essential to quantify the exact efficiency gap of SSL methods. Finally, to ensure these tools are safe for global deployment, future initiatives must prioritize large-scale, multi-center clinical validation to test model robustness across heterogeneous scanner protocols and diverse patient demographics, ultimately bridging the divide between academic benchmarks and reliable clinical practice.

6. Conclusion

The STS 2024 Challenge represents an advancement in semi-supervised instance segmentation for dental imaging, establishing the first standardized benchmark for label-efficient learning in this clinically critical domain. Our comprehensive analysis reveals that successful methodologies consistently leveraged synergistic combinations of three core strategies: knowledge transfer from foundation models, iterative pseudo-label refinement, consistency regularization learning, and multi-stage architectural optimization. The dominance of teacher-student frameworks employed by 85% of top-ranked solutions underscores their effectiveness in propagating knowledge from limited labeled data through carefully designed pseudo-labeling cycles. Notably, the integration of Segment Anything Model (SAM) variants and nnU-Net backbones emerged as particularly potent, enabling teams like Guo777 and Isjinhao to achieve instance-level Dice scores exceeding 90% in CBCT segmentation while utilizing merely 9% labeled data. These technical innovations directly translate to tangible clinical value, as evidenced by our human-machine collaboration study demonstrating that top solutions reduced annotation time by 88.2% for adult OPG and 94.1% for CBCT volumes when compared to manual annotation with junior clinicians achieving 73% of senior-expert efficiency when refining high-quality predictions.

Despite these advances, several critical challenges persist. The underutilization of unlabeled data by approximately 60% of participating teams highlights a persistent gap between semi-supervised learning potential and practical implementation. Furthermore, while multi-stage pipelines like Haoyuuu’s T3Net showed exceptional performance in 3D segmentation, their computational complexity presents deployment barriers in time-sensitive clinical settings. Limitations in dataset diversity, particularly the scarcity of pediatric CBCT scans and complex pathological cases, also constrain model generalizability. Looking forward, three priority directions emerge: 1) Developing hybrid SSL frameworks that combine SAM’s zero-shot generalization with uncertainty-aware consistency constraints, 2) Creating cross-modal learning architectures that leverage complementary information from both 2D and 3D dental imaging, and 3) Establishing standardized clinical validation protocols that assess both segmentation accuracy and real-world workflow integration. The publicly released STS 2024 dataset and open-sourced methodologies provide an essential foundation for addressing these priorities. By catalyzing research in label-efficient learning for dental instance segmentation, this initiative accelerates progress toward accessible, AI-enhanced diagnostics that can transform global oral healthcare, particularly in resource-constrained settings where expert annotation remains scarce.

CRedit authorship contribution statement

Yaqi Wang: Conceptualization, Investigation, Funding acquisition **Jun Liu:** Supervision, Project administration, Conceptualization, Funding acquisition, Writing – Review & Editing **Yifan Zhang:** Data Curation, Resources **Zhi Li:** Writing–Review & Editing, Validation, Visualization **Jiaxue Ni:** Investigation, Validation **Qian Luo:** Writing-Review, Validation **Jialuo Chen:** Writing-Review, Validation **Chengyu Wu:** Writing-Original Draft & Editing, Validation **Hongyuan Zhang:** Editing, Validation, Visualization **Jin Liu:** Investigation, Formal analysis **Can Han:** Investigation, Validation **Kaiwen Fu:** Competitor **Changkai Ji:** Competitor **Xinxu Cai:** Competitor **Jing Hao:** Competitor **Zhihao Zheng:** Competitor **Shi Xu:** Competitor **Junqiang Chen:** Competitor **Qianni Zhang:** Investigation, Formal analysis **Dahong Qian:** Investigation, Validation **Shuai Wang:** Supervision, Methodology, Software, Validation, Writing – Review & Editing **Huiyu Zhou:** Supervision, Project administration, Conceptualization, Review & Editing

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

This work was supported by the National Natural Science Foundation of China (No. 62206242, No. 62201323), Zhejiang Provincial Natural Science Foundation of China (No. LD25F020005), and China Science and Technology Foundation of Sichuan Province (No. 2022YFS0116). There are no conflicts of interest between authors. Yifan Zhang is the principal sponsor of the challenge by collecting and providing clinical data. Only the organizers and members of their immediate team have access to test case labels. The study protocol was approved by the Medical Ethics Committee of Hangzhou Stomatological Hospital (Approval No: 2022YR014).

Data availability

Data will be made available on request.

References

- [1] L. Sischo, H. Broder, Oral health-related quality of life: what, why, how, and future implications, *Journal of dental research* 90 (11) (2011) 1264–1270.
- [2] N. Shah, N. Bansal, A. Logani, Recent advances in imaging technologies in dentistry, *World journal of radiology* 6 (10) (2014) 794.
- [3] J. Cosson, Interpreting an orthopantomogram, *Australian Journal of General Practice* 49 (9) (2020) 550–555.
- [4] S. Jain, K. Choudhary, R. Nagi, S. Shukla, N. Kaur, D. Grover, New evolution of cone-beam computed tomography in dentistry: Combining digital technologies, *Imaging science in dentistry* 49 (3) (2019) 179.
- [5] X. Xu, C. Liu, Y. Zheng, 3d tooth segmentation and labeling using deep convolutional neural networks, *IEEE transactions on visualization and computer graphics* 25 (7) (2018) 2336–2348.
- [6] P. Lahoud, M. EzEldeen, T. Beznik, H. Willems, A. Leite, A. Van Gerven, R. Jacobs, Artificial intelligence for fast and accurate 3-dimensional tooth segmentation on cone-beam computed tomography, *Journal of Endodontics* 47 (5) (2021) 827–835.
- [7] R. Sapkota, A. Paudel, M. Karkee, Zero-shot automatic annotation and instance segmentation using llm-generated datasets: Eliminating field imaging and manual annotation for deep learning model development, *arXiv preprint arXiv:2411.11285* (2024).
- [8] J. Redmon, You only look once: Unified, real-time object detection, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [9] H. Qin, X. Zhan, Y. Li, Y. Zheng, Flexssl: A generic and efficient framework for semi-supervised learning, *arXiv preprint arXiv:2312.16892* (2023).
- [10] A. Kumarihami, S. Heshani, P. Sathyathas, R. Illeperuma, Development of brief image quality evaluation criteria for digital orthopantomography (opg) images in dental radiography, *Journal of Health Science* 6 (2018) 139–147.
- [11] T. Jiang, Y. Wang, S. Liu, Q. Zhang, L. Zhao, J. Sun, Instance recognition of street trees from urban point clouds using a three-stage neural network, *ISPRS Journal of Photogrammetry and Remote Sensing* 199 (2023) 305–334.

- [12] Y. Wang, Y. Zhang, X. Chen, S. Wang, D. Qian, F. Ye, F. Xu, H. Zhang, R. Dan, Q. Zhang, et al., Miccai 2023 sts challenge: A retrospective study of semi-supervised approaches for teeth segmentation, *Pattern Recognition* (2025) 112049.
- [13] A. Abdi, S. Kasaei, Panoramic dental x-rays with segmented mandibles, <https://doi.org/10.17632/hxt48yk462.1>, mendeley Data, V1 (2017). doi: 10.17632/hxt48yk462.1.
- [14] D. Budagam, A. Kumar, S. Ghosh, A. Shrivastav, A. Z. Imanbayev, I. R. Akhmetov, D. Kaplun, S. Antonov, A. Rychenkov, G. Cyganov, A. Sinitca, Instance segmentation and teeth classification in panoramic x-rays (2024). [arXiv:2406.03747](https://arxiv.org/abs/2406.03747).
- [15] K. Panetta, R. Rajendran, A. Ramesh, S. P. Rao, S. Agaian, Tufts dental database: a multimodal panoramic x-ray dataset for benchmarking diagnostic systems, *IEEE journal of biomedical and health informatics* 26 (4) (2021) 1650–1659.
- [16] J. C. M. Román, V. R. Fretes, C. G. Adorno, R. G. Silva, J. L. V. Noguera, H. Legal-Ayala, J. D. Mello-Román, R. D. E. Torres, J. Facon, Panoramic dental radiography image enhancement using multiscale mathematical morphology, *Sensors* 21 (9) (2021) 3110. doi:10.3390/s21093110.
- [17] I. E. Hamamci, S. Er, E. Simsar, A. E. Yuksel, S. Gultekin, S. D. Ozdemir, K. Yang, H. B. Li, S. Pati, B. Stadlinger, A. Mehl, M. Gundogar, B. Menze, Dentex: An abnormal tooth detection with dental enumeration and diagnosis benchmark for panoramic x-rays (2023). [arXiv:2305.19112](https://arxiv.org/abs/2305.19112).
URL <https://arxiv.org/abs/2305.19112>
- [18] A. W. Radio, vrad2 dataset, <https://universe.roboflow.com/arshs-workspace-radio/vrad2> (sep 2024).
URL <https://universe.roboflow.com/arshs-workspace-radio/vrad2>
- [19] Y. Chen, H. Du, Z. Yun, S. Yang, Z. Dai, L. Zhong, Q. Feng, W. Yang, Automatic segmentation of individual tooth in dental cbct images from tooth surface map by a multi-task fcn, *IEEE Access* 8 (2020) 97296–97309. doi: 10.1109/ACCESS.2020.2991799.
- [20] W. Cui, Y. Wang, Y. Li, D. Song, X. Zuo, J. Wang, Y. Zhang, H. Zhou, B. s. Chong, L. Zeng, et al., Ctooth+: A large-scale dental cone beam computed tomography dataset and benchmark for tooth volume segmentation, in: *MICCAI*

Workshop on Data Augmentation, Labelling, and Imperfections, Springer, 2022, pp. 64–73.

- [21] M. Cipriano, S. Allegretti, F. Bolelli, M. Di Bartolomeo, F. Pollastri, A. Pellacani, P. Minafra, A. Anesi, C. Grana, Deep segmentation of the mandibular canal: A new 3d annotated dataset of cbct volumes, *IEEE Access* 10 (2022) 11500–11510. doi:10.1109/ACCESS.2022.3144840.
- [22] Y. Huang, W. Liu, C. Yao, X. Miao, X. Guan, X. Lu, X. Liang, L. Ma, S. Tang, Z. Zhang, et al., A multimodal dental dataset facilitating machine learning research and clinic services, *Scientific Data* 11 (1) (2024) 1291.
- [23] L. Lumetti, V. Pipoli, F. Bolelli, E. Ficarra, C. Grana, Enhancing patch-based learning for the segmentation of the mandibular canal, *IEEE Access* 12 (2024) 79014–79024.
- [24] F. Bolelli, L. Lumetti, S. Vinayahalingam, M. Di Bartolomeo, A. Pellacani, K. Marchesini, N. Van Nistelrooij, P. Van Lierop, T. Xi, Y. Liu, et al., Segmenting the inferior alveolar canal in cbcts volumes: the toothfairy challenge, *IEEE Transactions on Medical Imaging* (2024).
- [25] F. Bolelli, K. Marchesini, N. van Nistelrooij, L. Lumetti, V. Pipoli, E. Ficarra, S. Vinayahalingam, C. Grana, Segmenting maxillofacial structures in cbct volumes, in: *Proceedings of the Computer Vision and Pattern Recognition Conference, 2025*, pp. 5238–5248.
- [26] L. T. Arsiwala-Scheppach, A. Chaurasia, A. Mueller, J. Krois, F. Schwendicke, Machine learning in dentistry: a scoping review, *Journal of Clinical Medicine* 12 (3) (2023) 937.
- [27] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, Springer, 2015, pp. 234–241.
- [28] F. Milletari, N. Navab, S.-A. Ahmadi, V-net: Fully convolutional neural networks for volumetric medical image segmentation, in: *2016 fourth international conference on 3D vision (3DV)*, Ieee, 2016, pp. 565–571.
- [29] F. Isensee, P. F. Jaeger, S. A. Kohl, J. Petersen, K. H. Maier-Hein, nnu-net: a self-configuring method for deep learning-based biomedical image segmentation, *Nature methods* 18 (2) (2021) 203–211.

- [30] G. Dot, A. Chaurasia, G. Dubois, C. Savoldelli, S. Haghighat, S. Azimian, A. R. Taramsari, G. Sivaramakrishnan, J. Issa, A. Dubey, et al., Dentalsegmentator: robust open source deep learning-based ct and cbct image segmentation, *Journal of Dentistry* 147 (2024) 105130.
- [31] C. Wang, J. Yang, B. Wu, R. Liu, P. Yu, Trans-vnet: Transformer-based tooth semantic segmentation in cbct images, *Biomedical Signal Processing and Control* 97 (2024) 106666.
- [32] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask r-cnn, in: *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [33] D.-H. Lee, et al., Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks, in: *Workshop on challenges in representation learning, ICML, Vol. 3, Atlanta*, 2013, p. 896.
- [34] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin, C.-L. Li, Fixmatch: Simplifying semi-supervised learning with consistency and confidence, *Advances in neural information processing systems* 33 (2020) 596–608.
- [35] A. Tarvainen, H. Valpola, Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results, *Advances in neural information processing systems* 30 (2017).
- [36] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, C. A. Raffel, Mixmatch: A holistic approach to semi-supervised learning, *Advances in neural information processing systems* 32 (2019).
- [37] T. Chen, S. Kornblith, M. Norouzi, G. Hinton, A simple framework for contrastive learning of visual representations, in: *International conference on machine learning, PmLR*, 2020, pp. 1597–1607.
- [38] X. Wang, S. Gao, K. Jiang, H. Zhang, L. Wang, F. Chen, J. Yu, F. Yang, Multi-level uncertainty aware learning for semi-supervised dental panoramic caries segmentation, *Neurocomputing* 540 (2023) 126208.
- [39] V. Cheplygina, M. De Bruijne, J. P. Pluim, Not-so-supervised: a survey of semi-supervised, multi-instance, and transfer learning in medical image analysis, *Medical image analysis* 54 (2019) 280–296.

- [40] N. Tajbakhsh, L. Jeyaseelan, Q. Li, J. N. Chiang, Z. Wu, X. Ding, Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation, *Medical image analysis* 63 (2020) 101693.
- [41] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, et al., Segment anything, in: *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 4015–4026.
- [42] J. Cheng, J. Ye, Z. Deng, J. Chen, T. Li, H. Wang, Y. Su, Z. Huang, J. Chen, L. Jiang, et al., Sam-med2d, *arXiv preprint arXiv:2308.16184* (2023).
- [43] X. Xu, J. Chen, J. Yin, Tooth instance segmentation and disease detection with uncertainty-aware contrastive learning and cross-scale attention, *IEEE Journal of Biomedical and Health Informatics* (2025) 1–12doi:10.1109/JBHI.2024.3525460.
- [44] L. Spector, Computer-aided dental implant planning, *Dental Clinics of North America* 52 (4) (2008) 761–775, vi. doi:10.1016/j.cden.2008.05.004.
- [45] K. Fu, C. Chang, J. Chen, Q. Hu, A self-training pipeline for semi-supervised 2d teeth instance segmentation, in: Y. Wang, D. Qian, S. Wang, A. Ben-Hamadou, S. Pujades, L. Lumetti, C. Grana, F. Bolelli (Eds.), *Supervised and Semi-supervised Multi-structure Segmentation and Landmark Detection in Dental Data*, Springer Nature Switzerland, Cham, 2025, pp. 156–165.
- [46] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [47] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, M. Wang, Swin-unet: Unet-like pure transformer for medical image segmentation, in: *European conference on computer vision*, Springer, 2022, pp. 205–218.