

# Transformer-based Tooth Alignment Prediction with Occlusion and Collision Constraints

Anonymous ICCV submission

Paper ID 7185

## Abstract

The planning of digital orthodontic treatment requires providing tooth alignment, which relays clinical experiences heavily and consumes a lot of time and labor to determine manually. In this work, we proposed an automatic tooth alignment neural network based on Swin-transformer. We first re-organized 3D point clouds based on dental arch lines and converted them into order-sorted multi-channel textures, improving both accuracy and efficiency. We then designed two new orthodontic loss functions that quantitatively evaluate the occlusal relationship between the upper and lower jaws. They are important clinical constraints, first introduced and lead to cutting-edge prediction accuracy. To train our network, we collected a large digital orthodontic dataset in more than 2 years, including various complex clinical cases. We will release this dataset after the paper's publication and believe it will benefit the community. Furthermore, we proposed two new orthodontic dataset augmentation methods considering tooth spatial distribution and occlusion. We compared our method with most SOTA methods using this dataset, and extensive ablation studies and experiments demonstrated the high accuracy and efficiency of our method.

## 1. Introduction

Tooth correction, medically known as orthodontics[32], primarily involves the use of metal braces[9] or clear aligners[18] to alleviate or rectify the conditions of dental misalignment and malformation. With the widespread adoption of digital acquisition technologies, computer-aided alignment design has been paid extensive attention, such as those based on intraoral scanners[35] and cone-beam computed tomography (CBCT)[2]. 3D tooth models are initially segmented individually[16, 24, 45, 46], and then repositioned by the clinician considering various alignment factors such as the extent of dental protrusion, dental skeletal relationship, and periodontal conditions of the patient, etc.

It heavily relies on the clinical expertise of orthodontists and is time-consuming, thereby significantly increasing the duration and cost of orthodontic treatment planning.

With the advancement of artificial intelligence, learning-based methods for tooth alignment are emerging rapidly[20], aiming at achieving fully automated tooth alignment. Among these methods, PointNet-based[33] ones are particularly representative[21, 23, 25, 44]. TANet[44] employs PointNet to construct a feature extraction module, encoding both jaw global information and teeth local information, and utilizes MLP to design regressors for predicting the position of each tooth. PSTN[23] utilizes both PointNet[33] and PointNet++[34] for feature encoding, refining features based on a combination of local and global latent vectors to regress tooth transformation parameters. TAalignNet[25], also based on PointNet encoders and MLP decoders, employs Squeeze-and-Excitation Blocks[13] and shared FC sequences for feature propagation to predict alignment parameters.

PointNet-based tooth alignment prediction methods showed great potential, but limitations have been revealed in representing local features of point clouds[40] recently. This paper introduces a more advanced shift window transformer (referred to as Swin-T). It incorporates sliding window operations and hierarchical merging design on the foundation of traditional vision transformers, addressing issues such as lower precision due to the variability of objects, excessive pixel count leading to high computational complexity, and low computational efficiency encountered in transformer-based methods. As teeth share similar sizes and structures, they are sampled into uniformly sized 3D point clouds and transformed into regular multi-channel data, forming ordered data. Building upon this data organization, this paper proposes a multi-level channel compression structure based on Swin-T (SWTBS) and an SWTP module to respectively extract global information of tooth centers and local information of tooth point clouds. Benefiting from the performance optimization of shift windows and communication between windows, features of individual teeth can mutually inform one another, gradually expand-

036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075

076 ing the receptive field and exhibiting excellent global control over entire dental arches[14], thereby achieving higher prediction accuracy.  
077

078  
079 We have collected a dataset of 855 orthodontic alignment plans in two years, to be released with the paper and source code. We also introduce constraint data augmentation during preprocessing, evaluated in our experiments. Our dataset, consisting of manually aligned initial tooth models, offers two advantages over previous methods using intraoral scan data: (1) the aligned labels correspond directly with the original scan data, avoiding complex point cloud matching, and (2) dentists' manual alignment is considered a more reasonable ground truth than post-orthodontic scan data.  
080

081 To summarize, our contributions are as follows:  
082

- 083 • A lightweight tooth alignment network based on Swin-T is designed to replace traditional 3D point cloud feature extraction encoders. It organizes scattered point clouds into regularly sized and orderly sorted multi-channel texture forms, ensuring high efficiency and seamless compatibility with complex scenarios such as missing teeth and wisdom teeth, surpassing the accuracy of the SOTA method in tooth alignment.  
084
- 085 • Two occlusal loss functions, the occlusal projecting overlap loss and occlusal distance uniformity loss, are designed based on medical domain knowledge. These functions enable more accurate and efficient quantitative measurement for the occlusal relationship between the upper and lower jaws.  
086
- 087 • An extensively annotated orthodontic alignment dataset, tailored to better suit the requirements of orthodontists, has been labeled. It will be released after the paper's publication and benefits the community. Additionally, two new orthodontic data augmentation methods considering tooth spatial distribution and occlusion are proposed to further increase the scale of training data.  
088

## 112 2. Related Works

### 113 2.1. Learning-based tooth alignment

114 Existing AI tooth alignment methods primarily use 3D point cloud data as input, rather than mesh or voxel data[27, 49].  
115 Early AI tooth alignment methods were mainly based on the PointNet[33] structure and its derivatives. TANet[44] utilizes PointNet[33] to encode the point cloud features of intraoral scan segmentation models, including global and local features. It then employs graph neural networks to connect and communicate tooth local features, regressing the 6DOF information[41] of teeth. PSTN[23] uses PointNet[33] to encode global and local features and PointNet++[34] to encode local features. After fusion, it uses a decoder designed based on PointNet[33] to regress orthodontic transformations of teeth. TAIGNet[25] achieves feature extrac-  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127

128 tion of 3D tooth models and tooth arrangement. It utilizes PointNet[33] as the feature encoder and employs fully connected layer sequences and SE blocks[13] for feature propagation, finally using fully connected layers to regress rotation and translation.  
129

130 Besides PointNet, recent AI tooth alignment methods also adopted emerging network structures such as DGCNN[43] and diffusion models[12]. Wang et al. proposed an improvement to TANet using tooth landmarks, where DGCNN was utilized to extract point cloud information[40]. They proposed a hierarchical regression using a three-layer graph neural network structure[36] to better predict the displacement transformation of each tooth, the landmark serves as a component of the tooth frame. TAPoseNet employed DGCNN to predict the local coordinate axes of teeth and utilized an autoencoder to extract geometric information[5]. Additionally, they proposed a multi-scale GCN[50] to characterize the spatial relationships between teeth at different levels, enabling more accurate prediction of the target positions of teeth. LETA[37] extracts features through latent encoding of dual branches (original data & ground truth), and predicts 6DoF of teeth by utilizing encoding differences. During training, GT point clouds are required, while only original point clouds are needed to complete prediction.  
131

132 Lei et al. employed probabilistic diffusion models[38] to iteratively denoise random variables, learning the distribution of transformation matrices for dental transitions from malocclusion to normal occlusion, thus achieving more realistic orthodontic predictions[21]. Furthermore, the network structure of TAIGNet mentioned above was actually proposed in image-based tooth alignment methods, iOrthoPredictor[25], which utilize three-dimensional geometric information encoded in the unsupervised generative model StyleGAN[15]. Through meaningful paths in latent space normals, alignment processes in image space are generated. Due to the complexity of the occlusal action between the upper and lower teeth, the calculation methods based on angles or center points have large errors and limited effectiveness, thus having deficiencies.  
133

### 167 2.2. Shift window transformer

168 Transformers[39] have achieved major advancements in NLP[29] and computer vision domains[19]. Vision Transformer (ViT)[6] directly applies self-attention to image patches, achieving strong results in classification without CNNs[17]. Swin-T[26] builds on ViT with movable windows, limiting sub-attention to non-overlapping local regions for greater efficiency. Swin3D[51] adapts Swin-T for 3D point clouds, converting sparse points into voxel grids and using farthest point sampling (FPS)[8, 28] and KNN pooling[22], though this is computationally intensive and misses advantages from serializing relative positions. Our  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178

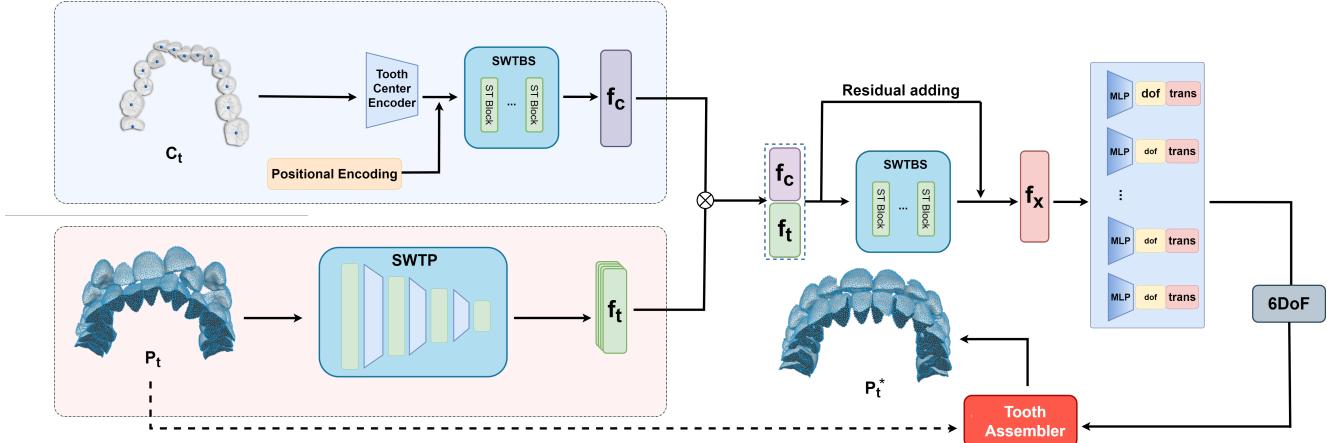


Figure 1. Network architecture overview. The encoding module has two branches: one for global features from the tooth center and one for local features from the tooth point cloud. Global features are extracted using SWTBS with shared Swin-T blocks, while local features are processed via SWTP with multi-stage hierarchical fusion. The features are merged, passed through SWTBS propagation, and then regressed by an MLP to predict the 6DOF transformation parameters for orthodontics

179 approach enhances the multi-level fusion network in Swin-  
 180 T[26] by using sliding windows to reduce data size layer by  
 181 layer, optimizing efficiency and expanding global receptive  
 182 fields. We use FPS[28] to uniformly sample and structure  
 183 each tooth’s data, avoiding direct KNN downsampling as  
 184 in[53] and[47], to improve information propagation.

### 3. Methodology

#### 3.1. Network overview

187 We segment the patient’s intraoral scan model to obtain the  
 188 gingival point cloud  $G$  and crown point clouds  $T$  for 32  
 189 teeth. These include up to 16 upper teeth (numbered 1 to 16)  
 190 and 16 lower teeth (numbered 17 to 32), following the uni-  
 191 versal tooth naming standard. Missing teeth are discussed  
 192 in Section 3.2. To ensure training efficiency, farthest point  
 193 sampling[30] is applied with  $N = 512$  to balance sampled  
 194 point quantity and network performance.

$$195 \quad T = \{t_i | 1 \leq i \leq 32\} \quad (1)$$

$$196 \quad P = \{p_j^t | t \in T, 1 \leq j \leq 512\} \quad (2)$$

198 Each tooth’s centroid is defined as  $C_t$ , where  $c_t$  is the  
 199 geometric centroid of the sampled point cloud  $P_t$ .

$$200 \quad C_t = \{c_t | t \in T, c_t = \text{Ave}(p)\}_{p \in P_t} \quad (3)$$

201 The goal of this study is to predict the 6Dof[7] pose  
 202 transformation parameters for each tooth model, using both  
 203 pre/post-orthodontic data. To ensure accurate loss calcu-  
 204 lation during training, the dataset includes corresponding  
 205 ground truth data  $P_t^*$ , with consistent sampling positions  
 206 and orders for  $P_t$  and  $P_t^*$ .

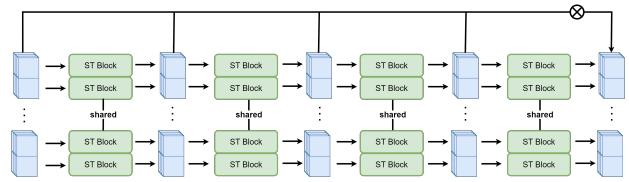


Figure 2. SWTBS module: Four groups of shared Swin-T blocks, each with 16 channels, with residuals added to the final output.

207 This paper proposes a dual-module architecture for fea-  
 208 ture extraction, as illustrated in Figure 1. The global module  
 209 encodes tooth center points using MLP layers and position-  
 210 al encoding, followed by Swin-T blocks (SWTBS) to produce  
 211  $f_c$ . The local module processes the 3D tooth point cloud  
 212 through the Swin transformer pipeline (SWTP), producing  
 213  $f_t$ . After pooling and merging  $f_t$  with  $f_c$ , the resulting high-  
 214 dimensional vector  $F = \{f_c, f_t\}$  is processed by SWTBS to  
 215 yield  $f_x$ , enhancing feature optimization. Finally, a down-  
 216 sampling regression module obtains the 6DoF transforma-  
 217 tion parameters for orthodontics.

218 The tooth point cloud is processed through hierarchi-  
 219 cal downsampling, similar to the Swin Transformer[26], in-  
 220 volving patch partitioning and feature merging across four  
 221 stages, as shown in Figure 3. Unlike the original network,  
 222 we keep the channel count constant to prevent excessive  
 223 feature dimension and loss during MLP conversion. Data  
 224 columns, rather than rows, are merged during feature pro-  
 225 cessing, as the first dimension represents the number of  
 226 teeth, and each tooth’s rotations and translations are unique.

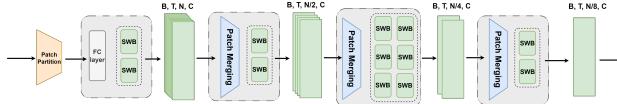


Figure 3. SWTP module, featuring a multi-stage feature fusion mechanism.

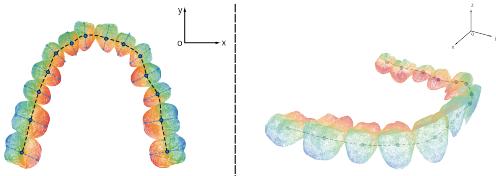


Figure 4. Points cloud are serialized according to their distance from simulated dental arch line, values on lingual side set as positive while labial side set as negative.

### 3.2. Data re-organization and augmentation

We designed a serialization method based on a simulated dental arch line, created by fitting central points from a tooth segmentation model and connecting them with Hermite curves. Points are sorted by their distance from the arch, with labial points positive and lingual points negative. Results in Figure 4 maintain relative positions among the 512 sampled points, improving network performance.

Regular data augmentation applies random rotation and translation on teeth based on a Gaussian distribution. It generates pre-orthodontic data while preserving the ground truth as post-orthodontic data. However, regular augmentation may produce clinical-illogical cases, including too far away from the arch lines and teeth collision. For this sake, we propose a constrained data augmentation that involves two relevant clinical constraints:

#### 3.2.1. Jaw regularization constraint

If the distance between two teeth exceeds 2.35 mm, the farther tooth is moved towards the central incisor along the dental arch line until the gap is within the threshold, as shown by the red teeth in Figure 5. Teeth exceeding 2.2 mm from the arch are pulled inward, as shown by the blue teeth, based on dataset statistics. Our strategy is to move the distal tooth towards the mesial tooth, specifically in the direction of the central incisor until the inter-tooth gap is within the threshold. This approach minimizes large gaps in the augmented dataset, ensuring more reasonable data. If the movement increases inter-tooth distance in the opposite direction or causes tooth collisions, the method in Section 3.2.2 is used for detection and avoidance, all along the simulated dental arch line.

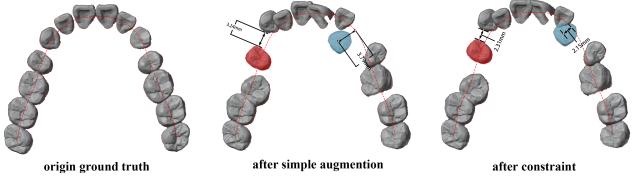


Figure 5. The maxillofacial regularization corrects excessive gaps or deviations based on dataset statistics.

#### 3.2.2. Collision detection constraints

We use a BVH collision detection algorithm[10, 31] to detect collisions and identify the colliding parts. Upon collision, the simulated dental arch line is used to avoid interlocking while preserving the arch shape.

For the efficiency purpose, we parallelize the BVH construction with a tooth-wise multi-threaded acceleration. It is worth noting that the BVH construction and collision detection are only employed in the pre-processing stage, which has more tolerance on the performance.

It is worth mentioning that missing teeth are specially handled, thus have very limited impact on both serialization and data augmentation. Though Missing teeth will affect the feature extraction in a single window, the impact is negligible for the sliding window that moves in the data.

### 3.3. Loss functions

The global loss function of the network consists of four components, the latter two are specifically designed in this paper to address dental occlusion. Each part of the loss function will be elaborated in the following part of this subsection.

$$L = \delta_0 * L_{recon} + \delta_1 * L_{fit} + \delta_2 * L_{uni} + \delta_3 * L_{val} \quad (4)$$

The hyperparameters  $\delta_0, \delta_1, \delta_2$  and  $\delta_3$  are used to weight each component accordingly.

#### 3.3.1. Reconstruction loss

We utilized the model reconstruction loss mentioned in[40]. Different from[40], the post-orthodontic data in our dataset was manually adjusted by orthodontists using the pre-orthodontic data. Therefore, the vertex positions and orders of the models before and after orthodontic treatment are correspond.

$$L_{recon}^{point} = \sum_{t \in T} \left( \sum_{i=0, p \in P_t} \|\bar{p}_i - p_i^*\|_2^2 + \|\bar{c}_t - c_t^*\|_2^2 \right) \quad (5)$$

#### 3.3.2. Transformation parameter loss

The transformation parameter loss comprises two components: rotation loss and translation loss. We computes

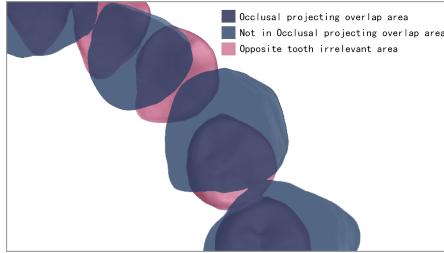


Figure 6. Visualization of the occlusion projection range.

weights for each tooth during training, based on the magnitude of misalignment. These weights are cumulatively added to the original loss, as shown in Equation 7 & 8. Drawing from [52], we emphasize that more severe misalignments should receive greater attention, corresponding to larger loss values.

$$L_{val} = \omega * L_{rotate} + L_{trans} \quad (6)$$

$$L_{rotate} = \sum_{t \in T} \left[ \sum_{i=0}^3 L_1 \left( \overline{rotate}_t(i), rotate_t^*(i) \right) * (1.0 + \zeta_t^{rotate}) \right] \quad (7)$$

$$L_{trans} = \sum_{t \in T} \left[ \sum_{i=0}^2 L_1 \left( \overline{trans}_t(i), trans_t^*(i) \right) * (1.0 + \zeta_t^{trans}) \right] \quad (8)$$

Due to the relatively small numerical values of rotation loss, an additional parameter  $w$  is introduced to amplify the impact of quaternion rotation loss during actual training.

### 3.3.3. Occlusal projecting overlap

Occlusion projection range consistency loss represents whether the interocclusal region between the predicted results of upper and lower jaws matches the ground truth. The definition of occlusion projection range is as follows: Let tooth  $t$  in one jaw have a corresponding area  $\beta_t$  in the opposite jaw. Project all points of  $t$  and  $\beta_t$  onto the occlusal plane.

$$m_i = \operatorname{Argmin}_{p_j \in P_{\beta_t}^f} \|p_i - p_j\|_2, p_i \in P_t^f \quad (9)$$

If the closest distance between points from  $t$ 's point cloud and  $\beta_t$ 's point cloud (on the occlusal plane) is less than a threshold  $\tau$ , then points from  $t$ 's point cloud within this distance are considered part of  $t$ 's occlusion projection range. We introduce the concept of occlusion projection range to bring predicted results closer to the ground truth at the occlusion projection range level. In the Figure 6,  $P_t^f$  represents the point cloud of tooth  $t$  projected onto the occlusal plane,  $P_{\beta_t}^f$  represents the point cloud of region  $\beta_t$  projected onto the occlusal plane, and  $m_i$  denotes the minimum

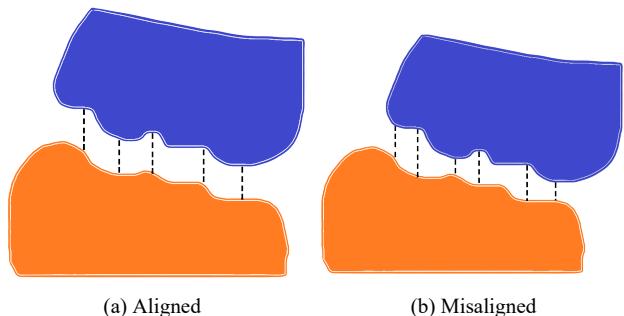


Figure 7. An illustration of occlusal distance uniformity. The variation of occlusion distances (black dotted lines) in aligned scenarios (a) is much smaller than the one in misaligned scenarios (b).

two-dimensional plane distance between a point  $p_i$  from the point cloud of  $t$  and the nearest point  $p_j$  from the point cloud of  $\beta_t$ .

$\tau$  is the threshold used to divide the occlusion projection range.  $X$  is a binary sequence where  $X_t(i)$  records whether point  $p_i$  from the point cloud of tooth  $t$  belongs to the occlusion projection range based on the relationship between  $m_i$  and  $\tau$ . If  $m_i$  is less than  $\tau$ , then  $X_t(i)$  takes the value of 1; otherwise, it takes the value of 0.

$$L_{fit} = \operatorname{Ave}_{t \in T} \left( \sum_{i=0}^{n-1} |\overline{X}_t(i) - X_t^*(i)| \right) \quad (10)$$

### 3.3.4. Occlusal distance uniformity

Due to the completely different morphologies and occlusion patterns between anterior and posterior teeth, we have proposed distinct loss function designs for anterior and posterior teeth based on discussions in [1, 4, 11]. Therefore, the occlusal distance uniformity loss  $L_{uni}$  across upper and lower jaws is composed of two parts:  $L_{uni}^{ant}$  for anterior teeth and  $L_{uni}^{post}$  for posterior teeth. As shown in the following formula, we introduce a weighting parameter  $w^{post}$  to balance.

$$L_{uni} = L_{uni}^{ant} + w^{post} \cdot L_{uni}^{post} \quad (11)$$

This paper evaluates the consistency and similarity of vectors connecting corresponding points in the occlusal regions of posterior teeth, based on the projection ambiguity constraint discussed in [4]. The occlusal distance uniformity is calculated within the occlusal projection range defined in Section 3.3.3. For a point  $p_i$  in the occlusal projection range of tooth  $t$ , the distance to the nearest point  $p_j$  in the corresponding range  $\beta_t$  on the opposing jaw is denoted as  $d$ , and the collection of all such distances forms set  $D$ . The uniformity of  $D$  represents the degree to which the occlusal ranges match concavely or convexly, defining the occlusal distance uniformity loss function for posterior teeth.

$$L_{uni}^{prior} = \sum_{t \in T_{prior}} Var_{X_t(i)=1} \left( \min_{X_{\beta_t}(j)=1} \|p_i - p_j\|_2 \right) \quad (12)$$

Due to the more prominent crowns of the upper incisors, their shape differs significantly from the molars, so the distance uniformity calculation method for molars cannot be applied. Based on the upper and lower anterior tooth correspondence described in [1, 11], we propose the vertex coordinate difference and angular difference between the tooth axis vector and the ground truth, denoted as  $L_{uni}^{ant1}$  and  $L_{uni}^{ant2}$ , respectively.

$$L_{uni}^{ant} = L_{uni}^{ant1} + \omega^{ant} * L_{uni}^{ant2} \quad (13)$$

$\overline{Peak}_t$  and  $Peak_t^*$  correspond to the highest incisal points in the predicted and ground truth data, respectively.

$$L_{uni}^{ant1} = \sum_{t \in T_{ant}} \|\bar{c}_t - c_t^*\|_2 + \sum_{t \in T_{ant}} \|\overline{Peak}_t - Peak_t^*\|_2 \quad (14)$$

$$L_{uni}^{ant2} = \sum_{t \in T_{ant}} \arccos \left( \frac{(\overline{Peak}_t - \bar{c}_t) \cdot (Peak_t^* - c_t^*)}{\|\overline{Peak}_t - \bar{c}_t\|_2 * \|Peak_t^* - c_t^*\|_2} \right) \quad (15)$$

## 4. Experiments

### 4.1. Dataset pre-processing and evaluation metrics

Our dataset contains 855 sets of dental data, which are derived from the 3D models of the upper and lower jaw teeth constructed through oral scans. Each set of dental data also includes data from multiple stages during the orthodontic treatment process, usually divided into about 30 stages. We can take the final orthodontic result of each patient as the ground tooth, and multiple treatment stages can respectively form multiple pairs of pre- and post-orthodontic data with the ground tooth. These data are first preliminarily segmented using the semantic segmentation network TSegNet[3], and then manually optimized for the mesh and edges. The optimized crown models are then aligned and arranged by experienced orthodontists to obtain the corresponding ground truth data. We randomly selected 700 samples for training, 35 samples for validation, and the remaining 120 samples for testing. Our labeled data do not come from the intraoral scan after orthodontic, because this would introduce differences in topology aspects on points before and after treatment.

We used an NVIDIA GeForce RTX 3090 (24GB VRAM) for 500-epoch training (batch size 8). Set N=512, initial learning rate 1.5e-4,  $w$  for the rotation loss in the transformation parameter loss was set to 10.0, The threshold  $\tau$  in Section 3.3.3 was empirically set to 0.07mm, shift window size  $8 \times 8$ . Our evaluation metrics used ADD/AUC from TANet[44] and landmark[40].

Table 1. Comparison of evaluation metrics between the proposed method and the SOTA method. Note that \* represents the effect of our method on the dataset[42].

Model	Test result			
	ADD ↓	ADD/AUC ↑	ME <sub>rotate</sub> ↓	ME <sub>translate</sub> ↓
TALigNet	1.5307	0.72	7.5461	2.0392
TANet	1.0075	0.81	6.9274	1.6815
PSTN	1.5889	0.71	8.6938	2.2155
Ptv3	1.2136	0.78	7.0663	1.7581
Landmark	0.8139	0.84	7.8277	1.3764
TADPM	1.1815	0.76	7.7426	1.7351
Ours*	<b>0.8115</b>	<b>0.84</b>	<b>2.9338</b>	<b>1.5904</b>
Ours	<b>0.6584</b>	<b>0.89</b>	<b>2.7678</b>	<b>1.1584</b>

## 4.2. Comparisons with SOTA methods

We tested the performance of some advanced methods on our dataset, including TANet[44], PSTN[23], TALigNet[25], Landmark[40], and TADPM[21]. Among them, the results of Landmark were obtained by the Wang et al. when they ran our dataset. We debugged the open-source code released by Lei et al. and reformatted our dataset according to the dataset format they published[42] to obtain the test results of TADPM. Although due to the limitations of the equipment, we used a more simplified dental mesh model and a smaller batch size, which may lead to some differences compared with the results published by Lei et al.[21], conducting the training and testing on the same hardware allows for a more rigorous comparison of the performance of all the methods. In addition, we processed their dataset into our format and conducted training and testing using our method, as shown in Table 1.

Training and testing data and specifications were consistent with Section 4.1. We compared ADD/AUC, average rotation error, and average translation error. Table 1 shows that our method performs best in all aspects, whether for AUC or rotational and translational deviations. Additionally, compared with the TADPM method, during the training and testing process, our method takes much less time to process a single case. Moreover, under the same time and equipment conditions, the quality of our method is better.

We compared the curves of average point distances, as shown in Figure 8. It is evident from the figure that our method achieves the highest accuracy under different definitions of average point distance. It is noteworthy that beyond an average point distance of 2.5, all curves converge to nearly 1.0. Therefore, the chart only displays curves for  $k \leq 2.5$ .

Figure 9 shows aligned tooth models achieved by our method versus others, with views from the front, side, and top. Our occlusal projection range alignment loss and occlusal distance uniformity loss help correct upper-lower jaw gaps and occlusal misalignment better than Chamfer vector loss. Landmarks improve focus on joint points, resolving misaligned gaps. Additionally, data serialization ensures ac-

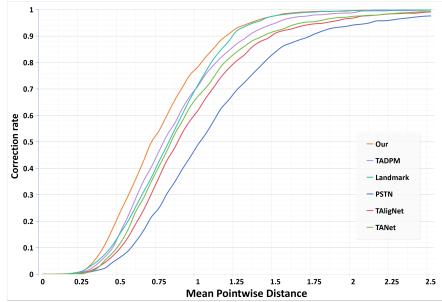


Figure 8. Comparison of accuracy curves between the proposed method and the SOTA method.

Table 2. Ablation experiment results of loss functions, testing the impact of different loss function combinations on final test results.

Loss fuc	Test result		
	$ADD/AUC \uparrow$	$ME_{rotate} \downarrow$	$ME_{translate} \downarrow$
$L_{recon}$	0.64	9.6	2.7
$L_{val}$	0.62	10.5	3.1
$L_{recon} + L_{val}$	0.79	8.3	2.2
$L_{recon} + L_{val} + L_{uni}$	0.81	5.9	1.7
$L_{recon} + L_{val} + L_{fit}$	0.83	5.3	1.4
$L_{recon} + L_{val} + L_{uni} + L_{fit}$	<b>0.89</b>	2.7	<b>1.1</b>

437 curate recognition of inter-tooth positions, even with incom-  
438 plete models, effectively handling issues with wisdom teeth  
439 and intra-jaw misalignment.

### 4.3. Visual results with complex cases

441 Figure 10 shows a comparison of our tooth alignment re-  
442 sults. It can be seen that our method produces very neat  
443 alignments even in complex situations such as large gaps,  
444 crossbite, and triangular tooth arrangements. The introduc-  
445 tion of serialization enhances the transformer’s ability to  
446 perceive the positions of teeth within the jaw, and the shift  
447 window efficiently extracts local features of the teeth, sup-  
448 ported by a comprehensive combination of loss functions.  
449 And our network specifying a maximum of 16 teeth per jaw,  
450 can handle cases with wisdom teeth or missing teeth.

### 4.4. Ablation study

#### 4.4.1. Loss functions

453 We discussed several loss functions in Section 3. Results  
454 of ablation experiments validating their effectiveness are  
455 as follows. Model reconstruction and transformation pa-  
456 rameter losses are effective. Our proposed occlusal losses  
457 improve prediction accuracy based on medical principles,  
458 which ensures that the upper and lower jaws are closely  
459 aligned according to natural occlusion laws and gradually  
460 move to their correct positions, despite longer training time.

#### 4.4.2. Network architecture

462 Table 3 shows that using the tooth center feature module, es-  
463 pecially SWTBS, yields better accuracy, as Swin-T’s sliding

Table 3. Ablation experiment results of network architecture.

Methods	w/o SWTP	$ADD/AUC \uparrow$	$ME_{rotate} \downarrow$	$ME_{translate} \downarrow$
VTBS	✓	0.79	7.10	1.83
PTv3	✓	/	/	/
SWTBS(Ours)	✓	<b>0.90</b>	<b>2.70</b>	<b>1.10</b>
VTBS	✗	0.75	8.50	2.20
PTv3	✗	0.73	9.80	2.40
SWTBS(Ours)	✗	0.81	7.20	2.00

Table 4. Ablation experiment results of serialization.

Serialization Function	Test result		
	$ADD/AUC \uparrow$	$ME_{rotate} \downarrow$	$ME_{translate} \downarrow$
Random Order	0.77	6.1	1.9
Based on dental local z-axis	0.80	5.4	1.7
Based on dental arch center	0.82	5.6	1.3
Based on virtual arch line	<b>0.89</b>	<b>2.7</b>	<b>1.1</b>

464 window and merging mechanisms enhance point cloud fu-  
465 nction. Without SWTP, accuracy decreases across modules,  
466 showing that SWTP’s multi-stage architecture better cap-  
467 tures global features. Using PTv3 alone further reduces per-  
468 formance due to lacking mechanisms like sliding windows,  
469 critical for effective inter-tooth feature extraction.

470 It can be observed that compared to network structures  
471 using only Swin blocks or Vision blocks, the multi-stage  
472 Swin block structure yields higher accuracy. This is because  
473 the multi-stage approach reduces the size of the latent vector  
474 progressively through dimensional merging, which is more  
475 effective in retaining task-specific dental features than di-  
476 rectly passing down through averaging.

#### 4.4.3. Point cloud serialization

477 We discussed sorting points of individual teeth for input into  
478 the Swin-T multi-layer feature fusion module in Section 3.  
479 Serialization ensures points selected by the window corre-  
480 sponds to the same local region of the teeth, can better ex-  
481 tract relative position features between teeth[48].

482 As shown in Table 4, the random sorting method was  
483 worst as it couldn’t use sequential information to boost the  
484 transformer’s performance. Sorting by the local Z-axis of  
485 the teeth had sequence benefits and better performance, but  
486 it was still not enough because of multi-peaked tooth crowns  
487 (common in posterior teeth), making sequentially arranged  
488 points in different local regions. The center-point-based  
489 sorting method solved this problem but had angular devia-  
490 tions for posterior teeth since the teeth are U-shaped. Our  
491 method based on the simulated dental arch curve for U-  
492 shaped jaws achieved best prediction results.

#### 4.4.4. Data augmentation

493 Our constrained data augmentation increases training the  
494 data scale, but excessive augmentation may lead to network  
495 distortion. To avoid it, we restrict the mount of augmented  
496 cases. Figure 11 shows the perdition accuracy with vary-  
497 ing augmentation ratio of augmented cases to original ones,  
498 it reaches the peak when the ratio is 54%. The accuray

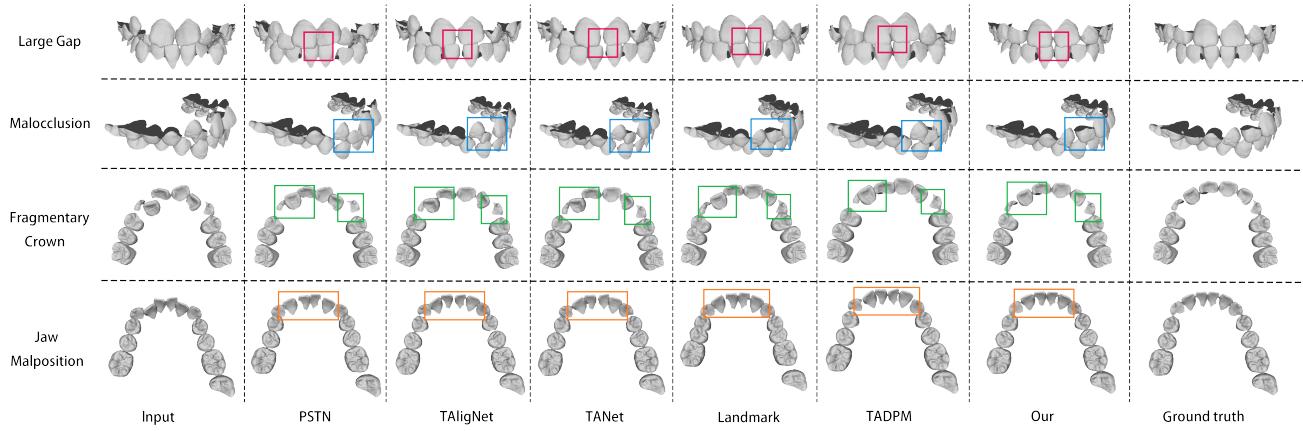


Figure 9. Comparison of prediction results between the proposed method and other methods. Here, four typical and challenging orthodontic problems are selected, listed from top to bottom: large gap, malocclusion, fragmentary crown, and jaw malposition.

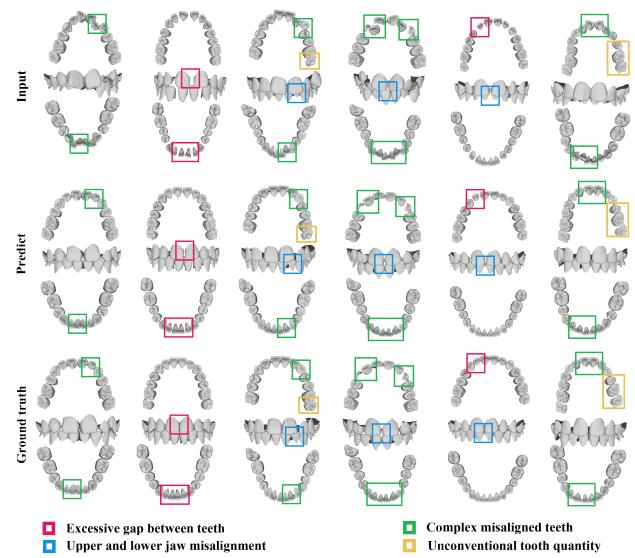


Figure 10. The figure shows the alignment prediction results of our method on 6 data cases in 6 columns.

of regular augmentation reaches the peak when the ratio is 62%. We compared our constrained augmentation with regular augmentation, on both source data and target data. Table 5 shows accuracy statistics with these scenarios, it gives two insights: (1) constraint augmentation contributes more significantly than regular augmentation, as regular augmentation ignores clinical requirements and introduces more biases, and (2) augmentation on target data provides better training quality than source data.

## 5. Conclusions

This paper proposes a novel, high-precision and efficient neural network approach for tooth alignment prediction. It

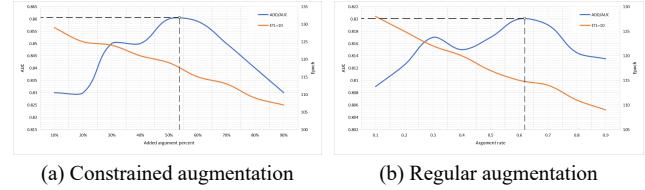


Figure 11. The effect of data augmentation intensity on final prediction accuracy (blue curves) and training convergence speed (orange curves).

Table 5. Ablation experiment results of data augmentation.

Augmentation	ADD/AUC ↑	ETL = 10 ↓	SigAugTime ↓
None	0.83	145	/
Source Data	0.84	125	0.84
	0.86	141	3.53
Target Data	0.85	109	0.84
	<b>0.90</b>	<b>134</b>	<b>3.51</b>

uses the multi-level feature fusion structure of Swin-T as its core, supplemented by a tooth center feature extraction module that emphasizes global features. Two occlusion evaluation loss functions are designed to effectively describe the occlusal relationships between upper and lower jaws. Furthermore, this paper constructed a open dataset in the field of tooth arrangement. This dataset includes over 855 fully annotated data pairs, consisting of point clouds sampled from tooth crowns, addressing the issue of a lack of public datasets in this field. A new constrained augmentation method is proposed to further augment the datasets. For future work, we plan to consider other stomatologic constraints for the tooth alignment task, and predict the full path of the tooth orthodontic treatment instead of the target position.

513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528

529 **References**

- [1] Lawrence F Andrews. The six keys to normal occlusion. *Am J orthod*, 62(3):296–309, 1972. 5, 6
- [2] Cheng Cheng, Xiaosheng Cheng, Ning Dai, Yi Liu, Qilei Fan, Yulin Hou, and Xiaotong Jiang. Personalized orthodontic accurate tooth arrangement system with complete teeth model. *Journal of medical systems*, 39:1–12, 2015. 1
- [3] Zhiming Cui, Changjian Li, Nenglun Chen, Guodong Wei, Runnan Chen, Yuanfeng Zhou, Dinggang Shen, and Wenping Wang. Tsegnet: An efficient and accurate tooth segmentation network on 3d dental model. *Medical Image Analysis*, 69: 101949, 2021. 6
- [4] S Davies and RMJ Gray. What is occlusion? *British dental journal*, 191(5):235–245, 2001. 5
- [5] Qingxin Deng, Xunyu Yang, Minghan Huang, Landu Jiang, and Dian Zhang. Taposenet: Teeth alignment based on pose estimation via multi-scale graph convolutional network. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 314–323. Springer, 2024. 2
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- [7] Guoguang Du, Kai Wang, Shiguo Lian, and Kaiyong Zhao. Vision-based robotic grasping from object localization, object pose estimation to grasp estimation for parallel grippers: a review. *Artificial Intelligence Review*, 54(3):1677–1734, 2021. 3
- [8] Yuval Eldar, Michael Lindenbaum, Moshe Porat, and Yehoshua Y Zeevi. The farthest point strategy for progressive image sampling. *IEEE transactions on image processing*, 6(9):1305–1315, 1997. 2
- [9] Theodore Eliades, Spiros Zinelis, Christoph Bourauel, and George Eliades. Manufacturing of orthodontic brackets: a review of metallurgical perspectives and applications. *Recent Patents on Materials Science*, 1(2):135–139, 2008. 1
- [10] Yan Gu, Yong He, Kayvon Fatahalian, and Guy Bleloch. Efficient bvh construction via approximate agglomerative clustering. In *Proceedings of the 5th High-Performance Graphics Conference*, pages 81–88, 2013. 4
- [11] LT Hiew, SH Ong, and Kelvin WC Foong. Optimal occlusion of teeth. In *2006 9th International Conference on Control, Automation, Robotics and Vision*, pages 1–5. IEEE, 2006. 5, 6
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2
- [13] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 1, 2
- [14] Silvana Allegrini Kairalla, Giuseppe Scuzzo, Tarcila Triviño, Leandro Velasco, Luca Lombardo, and Luiz Renato Paranhos. Determining shapes and dimensions of dental arches for the use of straight-wire arches in lingual technique. *Dental press journal of orthodontics*, 19:116–122, 2014. 2
- [15] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 2
- [16] Suzi Kim and Sunghee Choi. Automatic tooth segmentation of dental mesh using a transverse plane. In *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 4122–4125. IEEE, 2018. 1
- [17] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012. 2
- [18] Kislaya Kumar, Shivani Bhardwaj, and Vishal Garg. Invisalign: A transparent braces. *Journal of Advanced Medical and Dental Sciences Research*, 6(7):148–150, 2018. 1
- [19] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 2
- [20] Min Kyeong Lee, Veerasathpurush Allareddy, Sankeerth Rampa, Mohammed H Elnagar, Maysaa Oubaidin, Sumit Yadav, and Shankar Rengasamy. Applications and challenges of implementing artificial intelligence in orthodontics: A primer for orthodontists. In *Seminars in Orthodontics*. Elsevier, 2024. 1
- [21] Changsong Lei, Mengfei Xia, Shaofeng Wang, Yaqian Liang, Ran Yi, Yuhui Wen, and Yongjin Liu. Automatic tooth arrangement with joint features of point and mesh representations via diffusion probabilistic models. *arXiv preprint arXiv:2312.15139*, 2023. 1, 2, 6
- [22] Jiaxin Li, Ben M Chen, and Gim Hee Lee. So-net: Self-organizing network for point cloud analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9397–9406, 2018. 2
- [23] Xiaoshuang Li, Lei Bi, Jinman Kim, Tingyao Li, Peng Li, Ye Tian, Bin Sheng, and Dagan Feng. Malocclusion treatment planning via pointnet based spatial transformation network. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 105–114. Springer, 2020. 1, 2, 6
- [24] Sheng-hui Liao, Shi-jian Liu, Bei-ji Zou, Xi Ding, Ye Liang, and Jun-hui Huang. Automatic tooth segmentation of dental mesh based on harmonic fields. *BioMed research international*, 2015(1):187173, 2015. 1
- [25] YANG Lingchen, SHI Zefeng, Wu Yiqian, LI Xiang, ZHOU Kun, FU Hongbo, and Youyi Zheng. iorthopredictor: model-guided deep prediction of teeth alignment. *ACM Transactions on Graphics*, 39(6):216, 2020. 1, 2, 6
- [26] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 2, 3
- [27] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. 9

- 642 In 2015 IEEE/RSJ international conference on intelligent  
643 robots and systems (IROS), pages 922–928. IEEE, 2015. 2
- 644 [28] Facundo Mémoli and Guillermo Sapiro. A theoretical and  
645 computational framework for isometry invariant recognition  
646 of point cloud data. *Foundations of Computational Mathematics*,  
647 5:313–347, 2005. 2, 3
- 648 [29] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado,  
649 and Jeff Dean. Distributed representations of words and  
650 phrases and their compositionality. *Advances in neural in-*  
651 *formation processing systems*, 26, 2013. 2
- 652 [30] Carsten Moenning and Neil A Dodgson. Fast marching far-  
653thest point sampling. Technical report, University of Cam-  
654 bridge, Computer Laboratory, 2003. 3
- 655 [31] Jia Pan, Sachin Chitta, and Dinesh Manocha. Fcl: A general  
656 purpose library for collision and proximity queries. In 2012  
657 *IEEE International Conference on Robotics and Automation*,  
658 pages 3859–3866. IEEE, 2012. 4
- 659 [32] William R Proffit, Henry Fields, Brent Larson, and David M  
660 Sarver. *Contemporary Orthodontics, 6e: South Asia Edition-  
661 E-Book*. Elsevier Health Sciences, 2019. 1
- 662 [33] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas.  
663 Pointnet: Deep learning on point sets for 3d classification  
664 and segmentation. In *Proceedings of the IEEE conference on  
665 computer vision and pattern recognition*, pages 652–660,  
666 2017. 1, 2
- 667 [34] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J  
668 Guibas. Pointnet++: Deep hierarchical feature learning on  
669 point sets in a metric space. *Advances in neural information  
670 processing systems*, 30, 2017. 1, 2
- 671 [35] Marta Revilla-León, Dean E Kois, Jonathan M Zeitler, Wael  
672 Att, and John C Kois. An overview of the digital occlusion  
673 technologies: Intraoral scanners, jaw tracking systems, and  
674 computerized occlusal analysis devices. *Journal of Esthetic  
675 and Restorative Dentistry*, 35(5):735–744, 2023. 1
- 676 [36] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Ha-  
677 genbuchner, and Gabriele Monfardini. The graph neural net-  
678 work model. *IEEE transactions on neural networks*, 20(1):  
679 61–80, 2008. 2
- 680 [37] Zefeng Shi, Zijie Meng, Ruizhe Chen, Yang Feng, Zeyu  
681 Zhao, Jin Hao, Bing Fang, Zuozhu Liu, and Youyi Zheng.  
682 Leta: Tooth alignment prediction based on dual-branch latent  
683 encoding. *IEEE Transactions on Visualization and Computer  
684 Graphics*, 2024. 2
- 685 [38] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan,  
686 and Surya Ganguli. Deep unsupervised learning using  
687 nonequilibrium thermodynamics. In *International confer-  
688 ence on machine learning*, pages 2256–2265. pmlr, 2015. 2
- 689 [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko-  
690 reit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia  
691 Polosukhin. Attention is all you need. *Advances in neural  
692 information processing systems*, 30, 2017. 2
- 693 [40] Chen Wang, Guangshun Wei, Guodong Wei, Wenping Wang,  
694 and Yuanfeng Zhou. Tooth alignment network based on  
695 landmark constraints and hierarchical graph structure. *IEEE  
696 Transactions on Visualization and Computer Graphics*, 2022.  
697 1, 2, 4, 6
- 698 [41] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin,  
699 Shuran Song, and Leonidas J Guibas. Normalized object  
700 coordinate space for category-level 6d object pose and size  
701 estimation. In *Proceedings of the IEEE/CVF conference on  
702 computer vision and pattern recognition*, pages 2642–2651,  
703 2019. 2
- 704 [42] Shaofeng Wang, Changsong Lei, Yaqian Liang, Jun Sun, Xi-  
705 anju Xie, Yajie Wang, Feifei Zuo, Yuxin Bai, Song Li, and  
706 Yong-Jin Liu. A 3d dental model dataset with pre/post-  
707 orthodontic treatment for automatic tooth alignment. *Sci-  
708 entific Data*, 11(1):1277, 2024. 6
- 709 [43] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma,  
710 Michael M Bronstein, and Justin M Solomon. Dynamic  
711 graph cnn for learning on point clouds. *ACM Transactions  
712 on Graphics (tog)*, 38(5):1–12, 2019. 2
- 713 [44] Guodong Wei, Zhiming Cui, Yumeng Liu, Nenglun Chen,  
714 Runnan Chen, Guiqing Li, and Wenping Wang. Tanet:  
715 towards fully automatic tooth arrangement. In *Computer  
716 Vision–ECCV 2020: 16th European Conference, Glasgow,  
717 UK, August 23–28, 2020, Proceedings, Part XV 16*, pages  
718 481–497. Springer, 2020. 1, 2, 6
- 719 [45] Kan Wu, Li Chen, Jing Li, and Yanheng Zhou. Tooth  
720 segmentation on dental meshes using morphologic skeleton.  
721 *Computers & Graphics*, 38:199–211, 2014. 1
- 722 [46] Tai-Hsien Wu, Chunfeng Lian, Sanghee Lee, Matthew  
723 Pastewait, Christian Piers, Jie Liu, Fan Wang, Li Wang,  
724 Chiung-Ying Chiu, Wenchih Wang, et al. Two-stage mesh  
725 deep learning for automated tooth segmentation and land-  
726 mark localization on 3d intraoral scans. *IEEE transactions  
727 on medical imaging*, 41(11):3158–3166, 2022. 1
- 728 [47] Xiaoyang Wu, Yixing Lao, Li Jiang, Xihui Liu, and Heng-  
729 shuang Zhao. Point transformer v2: Grouped vector attention  
730 and partition-based pooling. *Advances in Neural Information  
731 Processing Systems*, 35:33330–33342, 2022. 3
- 732 [48] Xiaoyang Wu, Li Jiang, Peng-Shuai Wang, Zhijian Liu, Xi-  
733 hui Liu, Yu Qiao, Wanli Ouyang, Tong He, and Hengshuang  
734 Zhao. Point transformer v3: Simpler, faster, stronger. *arXiv  
735 preprint arXiv:2312.10035*, 2023. 7
- 736 [49] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Lin-  
737 guang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d  
738 shapenets: A deep representation for volumetric shapes. In  
739 *Proceedings of the IEEE conference on computer vision and  
740 pattern recognition*, pages 1912–1920, 2015. 2
- 741 [50] Zhilong Xiong and Jia Cai. Multi-scale graph con-  
742 volutional networks with self-attention. *arXiv preprint  
743 arXiv:2112.03262*, 2021. 2
- 744 [51] Yu-Qi Yang, Yu-Xiao Guo, Jian-Yu Xiong, Yang Liu, Hao  
745 Pan, Peng-Shuai Wang, Xin Tong, and Baining Guo. Swin3d:  
746 A pretrained transformer backbone for 3d indoor scene un-  
747 derstanding. *arXiv preprint arXiv:2304.06906*, 2023. 2
- 748 [52] yeyuxmf. Auto tooth arrangement. [https://github.com/  
749 yeyuxmf/auto\\_tooth\\_arrangement](https://github.com/yeyuxmf/auto_tooth_arrangement), 2022. 5
- 750 [53] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and  
751 Vladlen Koltun. Point transformer. In *Proceedings of  
752 the IEEE/CVF international conference on computer vision*,  
753 pages 16259–16268, 2021. 3