Anand Gupta, 2024, 12:3 ISSN (Online): 2348-4098 ISSN (Print): 2395-4752

An Open Access Journal

Shape-Guided Conditional Image Synthesis for Tooth Alignment

Anand Gupta¹, Deep Mehta², Arpan Dholakiya³, Soorajkumar Suresh⁴, Siddhant Jadhav⁵, Amit Yadav⁶, Yaduvir Singh⁷

Toothlens Healthcare Private limited, Delaware, United States^{1,2,3,4,5} 7Niet, Greater Noida, Uttar Pradesh, India^{1,6,7}

Abstract- Orthodontics is a field of dentistry that focuses on correcting misaligned teeth. Orthodontic treatment often involves the use of braces, aligners, retainers, and other dental appliances to align teeth, correct bite issues, and improve overall dental health and aesthetics. But fixing them can be a long and complicated process. So, it's important to have a picture showing what your teeth will look like after treatment. This helps dentists explain things better and makes patients more likely to go through with the treatment. In this paper, we present an approach for generating the conditional image of a misaligned tooth to an aligned tooth from a diffusion model it takes 2D images as input (e.g., captured by the camera, or smartphones), and aligns the teeth part in the image within 2D image space to generate post-treatment like images, our method employs 3D STL data scanned by intra-oral scanner projecting pre- and post-treatment, diffusion model to learn the spatial movement of each tooth followed by post-processing to convert the teeth area more realistic textures. We validate our pipeline over various facial photographs, illustrating its outstanding effectiveness and broad utility within the field of Orthodontics.

Keywords- Orthodontic Treatment, Pre- and Post-treatment Images, DDPM (Denoising Diffusion Probabilistic Models), Facial Recognition, Image Segmentation, smile design, Contour Segmentation, Edge detection, Conditional Diffusion Model, Tooth Alignment, Unet, Resnet18, Facial landmarks, Neural Network, Deep Learning, Image Synthesis, Shape-guided image synthesis, Loss Function, Training Pipeline, Inference, MaskGIT(Masked Generative Image Transformer), Data Pipeline, Image Generation, Visualization, Multilayer Perceptron(MLP)

I. INTRODUCTION

In recent years, there have been significant advances in Generative AI, particularly in the field of holds Image Generation, which promising implications for medical images. These advancements have led to the development of sophisticated models capable of generating realistic and high-fidelity images, thereby opening up new avenues for enhancing the quality and diversity of

medical images. As smile simulation gains increasing attention [3,4]. particularly in the realm of frontal face aesthetics, its potential to attract individuals to undergo orthodontic treatment becomes apparent. Patients opting for invisible teeth alignment procedures naturally desire insight into how their frontal facial appearance may evolve across various treatment stages. Orthodontic treatment offers an effective solution for correcting tooth misalignment and improving overall dental

© 2024 Anand Gupta. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (http://creativecommons.org/licenses/by/4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly credited.

alignment and aesthetics. while orthodontic intervention holds the potential to address malocclusion problems that affect over 90% of individuals to varying degrees, the benefits extend beyond mere physiological improvements. By enhancing oral health and function, orthodontic treatment plays a pivotal role in bolstering patients' confidence and overall psychological well-being. Nevertheless, the intricate and time-consuming nature of orthodontic procedures spanning months or even years may act as a deterrent for individuals considering treatments. Therefore, the creation and presentation of anticipated post-treatment facial images featuring aesthetically pleasing teeth assume paramount importance. These predictive visuals serve not only to captivate and inspire patients but also to facilitate more meaningful interactions between orthodontists and their patients.

While orthodontists may not be able to manually generate 2D teeth images based solely on their professional expertise and inform patients that these images represent the expected shapes of their teeth under the current treatment schedule, they can still offer reliable information. This includes providing 3D teeth models and their planning prescriptions, detailing how each tooth will move at every orthodontic stage. Therefore, the primary objective of our system is to reconstruct the teeth region of the facial images



Figure 1: we're showcasing before-and-after images of people's smiles. On the left we show the misaligned teeth and on the right post-treatment images of aligned teeth

In Traditional Clinical procedure orthodontists utilize a technique known as "Visual Treatment

Objective" (VTO) to visualize patients' appearance after orthodontic treatment, this method typically involves manipulating X-ray images to adjust soft tissues and skeleton based on specific landmarks [21,25]

the realm of computer vision, recent In advancements in deep learning, particularly in generative networks, have shown remarkable progress. However, many of these models rely heavily on paired images, which poses a challenge for our specific task. Collecting paired pre- and post-orthodontic facial photographs is arduous due to the lengthy orthodontic process and the evolving facial appearance over time. Moreover, 2D photographs lack the depth information necessary to represent the structure of teeth accurately. Consequently, bridging the gap between clinical knowledge of tooth alignment, which was defined on 3D tooth models, and 2D photographs presents a significant challenge [1]. In this research, our goal is to take photographs and generate aligned teeth and brighter teeth with realistic textures, as shown in Figure 1. Each generated photograph followed the unique tooth alignment property.

In this paper, we propose a shape-guided conditional image synthesis for tooth alignment in photographs. In essence, our main objective is to acquire the clinical understanding of tooth alignment as delineated in the 3D intra-oral scanning models [17], we followed the learned property to guide the post-orthodontic image generation. we use a pre-trained model which is trained by collecting a set of paired pre- and postorthodontic intra-oral scanning tooth models and applying them to the oral region of a 2D facial image. [1][22]. Post a Diffusion Model [11] approach is utilized to grasp information regarding tooth alignment. for generating tooth contours of the post-treatment image, threshold for Canny Edge (50,100) (as shown in Figure 2). Another diffusion model integrated to generate realistic is photographs with aligned teeth and incorporate teeth-lightening information. In this experiment, a substantial quantity of photographs is gathered from patients exhibiting varying degrees of

malocclusion issues. our method outperforms stateof-the-art techniques, including GANs [8].

II. LITERATURE REVIEW

Table 1: Summary of work done on generation of				
teeth alignment				

Author Name	Feature	Result and	Remark	
	Extraction	finding	D ¹ · ·	
Yulong Dou1, Lanzhuju Mei1, Dinggang Shen, Zhiming Cui	Faical landmarks, mouth detection, oral segmentatio n,tooth contours, skin color	State-of- the-art result for generating post treatment images	Pioneerin g approach in computer vision and conditiona l diffusion model to convert pre- treatment image to post- treatment image	
Richard P McLaughlin and John C Bennett[21]	Analyze orthodontic tooth movement, possibly introducing a new method	Dental VTO's accuracy in predicting tooth movements and its correlation	Practical Implicatio ns of the Dental VTO's and its potential impact on	
	called "dental	with actual treatment	clinical decision	
Shen Feihong, Liu jingjing, Li Haizhen[2]	The system utilizes image- based facial editing techniques to simulate orthodontic treatment effects accurately	outcomesEffectivenessoftheproposedapproach inpredictingtreatmentoutcomes,withOrthoGANoutperformingothermodelsintermsofsimulationauthenticity	Deep learning- based smile simulation approach for virtual teeth alignment effects in frontal face images with high precision generate realistic output	E E J J L H E E E E E E E E E E E E E E E E E E
Lingchen Yang, Zefeng Shi, Yiqian	Teeth geometry, enabling	iOrthoPredic tor accurately	iOrthoPred ictor offers a	 (]

Wu, Xiang Li,	fine-scale	predicts	significant
Kun Zhou,	geometry	teeth	advancem
Hongbo Fu,	control and	alignment,	ent in
and Youyi	lightening	demonstrati	predicting
Zheng [3]	condition.	ng practical	teeth
		applicability	alignment
		orthodontics	TSynNot
		through	TAliaNet T
		validation	GeoNet
		and user	
		studies	
Beijia Chen,	Teeth region	Demonstrat	StyleGAN
Hongbo Fu,	of portrait	es alignment	and
Kun Zhou,	photograph	prediction	BlendingN
and Youyi	s and	across	et enhance
Zheng[4]	embeds	diverse	alignment
	the latent	cases and	and
	space of	short video	blending
	StyleGAN.	clips	for
	facilitating,g		versatile
	eometrically		application
	_		•
Zhiming Cui,	Utilizing a	Represented	Tooth
Changjian Li,	two-stage	3D teeth	segmentat
and Wenping	network:	models	ion and
Wang [5]	edge	reconstructe	identificati
	and 2D	images	
	region	along with	images
	proposal.	augmenting	inages
	h h	the base	
		network	
Zhiming Cui,	ROI	Using deep	Advances
Yu Fang,	Extraction,	learning	in tooth
Lanzhuju Mei,	Tooth ROI,	achieve high	and bone
Bojun Zhang,	Tooth	accuracy	segmentat
BO YU,	centerea,	fully	ion
Caiwen liang	Skeleton	tooth and	
Vuhang Sun	Tooth	alveolar	
Lei Ma, liawei	classification	bone	
Huang, et	, Multi-task	segmentatio	
al[5]	tooth	n	
	segmentatio		
	n		
W De Vos, Jan	Reviewed	Identified	Suggested
Casselman,Sw	literature on	clinical	guidelines
ennen [7]	CBCI	applications	tor CBCT
	imaging	ana	device
	techniques	the lack of	yarameter s
		evidence-	ي.
		based data	
		on radiation	
		dose.	
lan	The	GAN	Compatibl
Goodfellow,	generative	training	e with
lean Pouget-	model	improves	various

Abadie,capturesMehdi Mirza,dataBing Xu,distributiDavid Warde-andFarley,SherjilOzair,AaroneCourville,andYoshuatheBengio[8]probabiloftrai	the model and on discriminato the r (The ativ adversarial odel process s refines both models) ty	image generation framework s for generative model	YANG Lingchen, SHI Zefeng, WU Yiqian, LI Viaga ZUQU	Frontal face image, misalignmen t teeth	model achieves high-quality image generation iOrthoPredic tor accurately predicts	ssion scheme The network has the potential for actiont
data Andreas Photo- Wirtz, Florian specific Jung, coupled Matthias Noll, shape Anqi Wang, models and Stefan model Wesarg[9] extracts tooth contours	2D Promising 2D 3D tooth reconstructi on with 0.848mm average symmetric surface distance, the	Automatic 3D reconstruc tion from orthodonti c photograp hs is feasible,Th	Kun, FU Hongbo, and Youyi Zhen[12]	with teeth geometry	alignment in photos, Disentangle ment of teeth geometry and in- mouth appearance is key	consultatio ns and treatment planning from the post- treatment image generation
from a image	each method achieves accurate reconstructi ons	is approach offers the potential for automatin g orthodonti c analysis.	Dor Bank, Noam Koenigstein, Raja Giryes[13]	Learn from edges, shapes, textures, colors, spatial relationships , anomalies	Effectively learn compressed and meaningful data representati ons	Challenges remain in setting hidden state size and distributio n for generative
Chengiei Wu, Parametr Derek tooth Bradley, Pablo prior, Garrido, model Michael a Zollhöfer, learned Christian database Theobalt, tooth Markus H shapes	row invasive the reconstructi uses on of entire pre- tooth rows from mouth photos, the method to creates	ne first approach for photo- based tooth row reconstruc tion, applicable to various	Diederik P. Kingma, Max Welling[14]	Extract latent features from images	Efficient and scalable gradient- based variational posterior inference	models VAE remains a widely used tool for probabilist ic modeling
Gross, and guide Thabo reconstru Beele[10] on	realistic icti tooth models from simple images	image sources, this innovative method offers a less invasive alternative to dental scanners for creating tooth models	Guillaume Alain, Yoshua Bengio, Li Yao, Jason Yosinski, Eric Thibodeau- Laufer, Saizheng Zhang, and Pascal Vincent[15]	GSN captures underlying data structures, facilitating tasks like denoising and generative modeling	GSN simplifies generative model training, capturing data distribution via Markov chain. Training sufficient- capacity models to denoise and	GSN extends beyond auto- encoders, aids classificati on, and offers fast sampling in deep models
Jonathan Ho, Appeara Ajay Jain, and geometr Pieter high-qua Abbeel[11] image synthesis	nce, State-of- y, the-art lity image synthesis on CIFAR10 and	Promising diffusion models with progressiv			reconstruct suffices for distribution capture	

Chituran	Litilizor	Unified	Dromotor
Chitwan	Utilizes	Unified	Promotes
Saharia,	conditional	framework	standardiz
William Chan,	diffusion	based on	ed
Huiwen	models for	conditional	evaluation
Chang, Chris	image-to-	diffusion	in image-
Lee, Jonathan	image	models	to-image
Ho, Tim	translation,	outperforms	translation;
Salimans,	surpassing	GAN and	generalist
David Fleet.	baselines	rearession	diffusion
and	without	haselines on	model
Mohammad	task-specific	image-to-	matches
Norouri[16]	tuning	image-to-	charcialist
Norouzi[10]	turning	inage translation	specialist
		translation	models
		tasks	performan
			ce.
Hidemichi	Evaluation	Intraoral	Intraoral
Kihara,	of accuracy	scanners'	scanner
Wataru	and	accuracy is	accuracy
Hatakeyama,	practicality	influenced	depends
Futoshi	of intraoral	bv lighting.	on the
Komine.	scanners	Repeatable	lighting.
Kvoko	and	for partial	Good for
Takafuii	verification	prostheses	partial
Toshiyuki	methods	Challenges	adantulau
Tokabachi lun	methous.	for full arch	edentation
Takanashi, Jun		ior iuli-arch	s cases
Yokota, Kenta		prostheses.	
Oriso, and			
Hisatomo			
Kondo[17]			
Vahid Kazemi	Ensemble of	Achieved	Utilizing
and Josephine	regression	super-	gradient
Sullivan[18]	trees	realtime	boosting
	estimates	performance	framework
	landmark	with high-	with
	positions	quality	appropriat
	from sparse	predictions.	e priors for
	nixel	handling	efficient
		nananng	cificient
	intensities	missing or	feature
	intensities	missing or	feature
	intensities for face	missing or partially	feature selection
	intensities for face alignment	missing or partially labeled data.	feature selection and
	intensities for face alignment	missing or partially labeled data.	feature selection and combating
	intensities for face alignment	missing or partially labeled data.	feature selection and combating overfitting.
Changqian	intensities for face alignment Bilateral	missing or partially labeled data. BiSeNet	feature selection and combating overfitting. BiSeNet
Changqian Yu, Jingbo	intensities for face alignment Bilateral Segmentatio	missing or partially labeled data. BiSeNet achieves	feature selection and combating overfitting. BiSeNet offers a
Changqian Yu, Jingbo Wang, Chao	intensities for face alignment Bilateral Segmentatio n Network	missing or partially labeled data. BiSeNet achieves 68.4% Mean	feature selection and combating overfitting. BiSeNet offers a promising
Changqian Yu, Jingbo Wang, Chao Peng,	intensities for face alignment Bilateral Segmentatio n Network balances	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on	feature selection and combating overfitting. BiSeNet offers a promising solution to
Changqian Yu, Jingbo Wang, Chao Peng, Changxin	intensities for face alignment Bilateral Segmentatio n Network balances spatial	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade-
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu,	intensities for face alignment Bilateral Segmentatio n Network balances spatial information	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS,	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20]	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20]	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20]	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20]	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic segmentatio	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion accuracy
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20]	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic segmentatio n	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion accuracy
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20]	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic segmentatio n.	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion accuracy
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20] Phillip Isola, Iun Yan Zhu	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic segmentatio n. Conditional adversarial	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion accuracy Shift from band
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20] Phillip Isola, Jun-Yan Zhu, Tingbui Zhau	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic segmentatio n. Conditional adversarial	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance Effective synthesis of	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion accuracy Shift from hand-
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou,	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic segmentatio n. Conditional adversarial networks	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance Effective synthesis of photos from label.	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion accuracy Shift from hand- engineere
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic segmentatio n. Conditional adversarial networks learn	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance Effective synthesis of photos from label maps,	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion accuracy Shift from hand- engineere d mapping
Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang[20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros[29]	intensities for face alignment Bilateral Segmentatio n Network balances spatial information and receptive field for semantic segmentatio n. Conditional adversarial networks learn mapping	missing or partially labeled data. BiSeNet achieves 68.4% Mean IOU on Cityscapes with 105 FPS, balancing speed and performance Effective synthesis of photos from label maps, object	feature selection and combating overfitting. BiSeNet offers a promising solution to the trade- off between speed and segmentat ion accuracy Shift from hand- engineere d mapping and loss

	functions, versatile for image-to- image translation tasks	on from edge maps, and image colorization	to automated learning in image translation
Vedant Singh, Surgan Jandial, Ayush Chopra, Siddharth Ramesh, Balaji Krishnamurth y, and Vineeth N Balasubraman ian[37]	Conditional diffusion models for image editing, stock photo generation, and 3-D object creation.	Conditional noise input enhances image control, diverging from traditional Gaussian noise approaches, input improves semantic attribute conditioning in diffusion models for image generation	Promising approach enhances image control and conditioni ng in diffusion models for diverse image generation
Olaf Ronneberger, Philipp Fischer, and Thomas Brox[51]	Network combines contracting and expanding paths for context and localization	Outperforms sliding- window network on ISBI challenge for neuronal structure segmentatio n	Fast segmentat ion, trained from few images, wins ISBI cell tracking challenge
Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Muller, Harry Saini Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, Robin Rombach[69]	Improved noise sampling techniques bias towards perceptually relevant scales in rectified flow models	Superior performance in high- resolution text-to- image synthesis compared to established diffusion formulations	Transform er-based architectur e enables bidirection al flow between image and text tokens, improving synthesis and comprehe nsion

In recent years, there has been a strong push towards digitizing the dental field. Intraoral scanners, which create impressions digitally, are seen as transformative for future dental treatments [17]. CBCT [7] provides dentists about information, the structure of patient's teeth. Various novel methods have emerged in associated domains, such as tooth segmentation, 3D tooth reconstruction [9],[10] and 3D tooth arrangement [35]. Lingchen et al [12] present a novel approach to visually predict teeth alignment in photographs along with a corresponding 3D teeth model as input and generate using the concept of StyleGAN. Feihong et al. [2] conclude their study with a highprecision visualization for orthodontic smile simulation, the system utilizes advanced imagebased facial editing techniques to provide highly realistic simulation results. Experiments and user studies confirm its effectiveness in predicting treatment outcomes in digital orthodontics.

Image generation in computer vision involves creating new digital images using algorithms or models, such as GANs [8]. Bank et al. [13] Autoencoder design to compress input data and then reconstruct it to resemble the original input. VAE [14], Denoising Diffusion Probabilistic Models [11] Saharia et al. [16] present image-to-image translation based on conditional diffusion models.

III. METHODOLOGY

1. Overview

In this research paper, we embark on a comprehensive exploration of Teeth alignment reconstruction. Our goal is to develop a tooth alignment multilayer network that understands the spatial movement of each instance of teeth that guides the generation of post-treatment photographs.

We have pre and post-intra-oral scanning pretrained models $S = \{S1, S2, ..., SN\}$ of the same patients. And facial photographs $I = \{I1, I2, ..., IM\}$ captured by smartphones. We used a module that is design for tooth alignment which seamlessly integrates 3D structural insights from intra-oral scanning models. This module dynamically selects

pre- and post-orthodontic intra-oral scanning models (i.e., $Sr \in S$ and $S^{\circ} r \in S^{\circ}$) for a given unpaired facial photograph (Ir \in I). It then executes a coarse 2D-3D registration. Then the 3D structure is projected onto the 2D facial photograph to obtain the contours of pre- and post-orthodontic $Cr \in R$ 128×256.

Post we used a pre-trained segmentation module to locate the mouth region and segment the binary mask and tooth contours from frontal face photographs. In addition, a generation model used which is from the conditional self-attention diffusion model along with Unet architecture to the realistic images After the initial processing, we apply a post-processing technique using an instance segmentation model to focus on specific instances and extract individual tooth masks from the segmented masks. Subsequently, we determine the shade of each tooth by identifying individual teeth and computing the mean saturation of the masked regions. We then adjust the H (Hue) and S (Saturation) channels to shift the tooth shade grade accordingly. The threshold list used for this adjustment includes values of 20, 40, 50, and 65.

Yulong Dou1 et al.



Figure 2: Flow of pipeline. When a 2D photograph is input into networks. It goes to the mouth detection module to obtain the mouth mask inside the lip area, and tooth contours. Then it passes through Align-Mod to predict we aligned tooth contours Lastly, the processed image undergoes the Generation Module (Gen-Mod), which generates a facial photograph displaying wellaligned teeth[1].

2. Image Segmentation

Detect face alignment (i.e position of the face) [18][19] and get the face size of Fi \in R 512×512

from input facial image Ii \in I \subseteq R. and locate the inside lip area from mouth segmentation model trained on ResNet18 backbone OD(.), YOLO v8. we propose a mouth detection model [20][23], segment the mask, and crop inside the lip area mask Ri \in R 128×256, for segmentation of the tooth we employ U-Net[51] Attention Refinement Module(ARMs), calculation of segmentation is as

$$Ri, Mi = OD(Ii),$$

$$Ci = U (Ri, Mi), \forall i = 1, 2,...,N,$$
 (1)

Where OD(.) represents mouth detection in photographs li, U(•) represent U-Net network. we obtain the pre-orthodontic tooth contours (Ci), mouth region (Ri), and oral mask (Mi) using the mouth segmentation model.

We used pre-trained network training utilizing a combination of Dice Loss [22] and Weighted Cross-Entropy Loss [24] to handle class imbalance effectively. Dice Loss emphasizes foreground learning, ideal for imbalanced scenarios, while Weighted Cross-Entropy Loss adjusts weighted proportions to address imbalance concerns. The loss function:

3. Orthodontic Alignment

We refer to the approach that innovatively To generate a complete mouth region, we integrates 3D intra-oral scan data into the Align-Mod model [1], crucial for orthodontic treatment planning. Unlike prior methods focused on reconstructing 3D dental models from tooth photographs, we project intra-oral scans onto 2D facial images using a 3D-to-2D rendering approach. Additionally, we've devised a conditional DDPM Model-based network to glean orthodontic insights from rendered tooth contours (Cr) achieved by applying edge detection that detects strong gradients with the threshold of 50-100 [9][27].

We apply coarse registration for establishing spatial correspondences between photographs based on the tooth landmarks [28]. Therefore, the standard contours Cr, C[^] r ∈ R 128×256

$$\rho m = (KR^T | -KR^TC) (M 1),$$
 (3)

where p is projection, K denotes the matrix of internal parameters, R denotes the rotation matrix indicating the camera's orientation, C is the camera spatial location M denotes the coordinates of the point in the world coordinate system, symbolized as $M = (X, Y, Z)^{\wedge} T$ and m denotes the coordinate of corresponding point in pixel coordinate system as $m = (u, v, 1)^T$.

4. Image-to-Image Generation

We refer image to image DDPM model [16] for the generation of the mouth regions with Gaussian noise Gr \in R^A 128×256 and we fine-tune the parameter beta schedule the number of time steps is 150 in mask-to-mask generation, contour-tocontour generation combined with conditional diffusion model. During training post orthodontic tooth contours $C^r = T$ (Cr c Gr), where c represents concatenation(channel-wise), Т represents learn the clinical knowledge.

Contours derived from the Image segmentation model (Ci)with Gaussian noise (Gi) applied into a diffusion model anticipating a dependable prediction for straight well-aligned teeth contours (C^i).

concatenate conditional information skin color Ki ∈ R 128×256, straight tooth contour, and Gaussian noise Gi ∈ R 128×256 to feed in our generation model.

5. Enhance Teeth

We proposed a teeth color enhancement network to focus on the teeth area. trained instance segmentation (Is) model, segmentation results are post-processed to obtain masks for individual teeth. Cropped mask from generated image and calculated the mean saturation of the masked image to identify tooth shade. The formula for mean saturation is:

$$M^{s} = \Sigma(Y^{s}) / N$$
 (4)

Where M^s represent mean saturation, Y^s is 1. Comparison yellow saturation filtered out zero values from the saturation channel of the HSV representation of the input image, and N is the number of filtered saturation values.

We then adjusted the tooth shade based on a threshold. modify the Hue (H) and Saturation (S) channels of the HSV image to shift the tooth shade grade. If the pixel was part of the tooth (determined by the mask), the saturation was updated according to the required shade. Additionally, all non-yellow hues were converted to yellow hues to handle edge cases. The modified HSV image was then converted back to the BGR color space to obtain the final result.

IV. EXPERIMENTS

we use a pre-trained model for facial segmentation to the crop mouth area, for mouth segmentation model and generation model [1], and trained mouth detection and semantic segmentation model for segmenting teeth used in postprocessing that generate more realistic textures. we set the batch size of 60.

V. RESULT



(a) mouth crop (b)mouth region (c)contours (d)straight tooth (c)generated teeth (f)changed shade Figure3: Aligned teeth generation process, crop the mouth(a), segment mouth region(b), tooth edges(c), generated tooth edges (d), generated mouth region(e), adjusted shade of teeth

We were able to generate the post-treatment image from 2d image captured by smartphones from our experiment with conditional image-toimage generation method using a diffusion model. The results obtained from our experiment are not just promising, demonstrating the potential of our approach. This success paves the way for further advancements in this field.

We conducted a comparative analysis of our conditional image synthesis tooth alignment model with Pix2Pix GAN [8][29], StyleGAN [3], and 3D Structure-guided Network [1]. The results revealed that our conditional diffusion post-processing network produces more realistic outputs compared to Pix2Pix, StyleGAN, and the 3D Structure-guided Network.

2. Discussion

In this experiment, we propose a shape-guided conditional image synthesis pipeline to generate post-treatment images from pre-treatment 2d photographs. As per experimental results give clinical significance. Additionally, we introduce preprocessing and post-processing approach to feed in Diffusion Model in task of mouth region image generation.

We trained our own instance segmentation model for individual teeth detection to detect and classify individual teeth for generating more shiny images. approach preserves the potential of Our orthodontic knowledge and ensures ease of use in practice within the orthodontic industry.

The limitation of pre-trained models is evident in our approach, as it struggles with severely misaligned cases and contour creation. While it successfully resolves merging issues in (collision and occlusion) in some cases, it remains unresolved for others. However, our method effectively addresses the generation of realistic shiny teeth colors and ensures that gums are generated with natural skin tones.

There are some guidelines for input images due to of the dlib our use shape_predictor_68_face_landmarks detection pretrained model, it will be unable to generate if the frontal face is missing and unable to generate from side views image. Additionally, if two or more consecutive teeth are missing, the model may not be able to generate these individual teeth accurately. These limitations highlight areas for future research and improvement

VI. FUTURE SCOPE

The future scope for the Shape-Guided Conditional Image Synthesis using the diffusion model is incredibly promising in the healthcare industry and spans various domains. Here are several domains where Shape-Guided Conditional Image Synthesis using a diffusion model is expected to have a significant impact.

Post-Treatment Reconstruction: Skin Diseases and Tumor Residuals: Shape-Guided Conditional 2. Synthesis image generation has potential to reconstruct the pre-treatment image of Skin Diseases and Tumor Residuals.

Cosmetic Surgery: This approach could be used to simulate the result of cosmetic procedures that can help patients make informed decisions about their treatment.

Virtual Reality: Image synthesis generation can contribute to developing realistic virtual environment for enlarging the shape of objects.

Fashion Design: This can be used for simulating the 5. mask of the clothes using shape and size and design. It could generate realistic images of how different clothing items or masks would fit.

VII. CONCLUSION

In conclusion, this project has presented comprehensive research and implementation of Shape Guided Image-Image to generation using conditional diffusion models. We started by exploring the fundamental concepts of smile design principles, and teeth alignment patterns through deep learning and moved to the generation from diffusion models. The process encompassed from data collection, preprocessing of data, facial landmarks detection, mouth detection, mouth segmentation, teeth instance segmentation, imageto-image generation, and postprocessing. our approach stands out from existing methods this ensures the reliability of our approach and its potential applicability in the healthcare sector.

Acknowledgements

This study was supported by the Noida Institute of Engineering and Technology (Department of Artificial Intelligence) and Toothlens Healthcare Private Limited.

REFERENCES

- 1. Yulong Dou1, Lanzhuju Mei1, Dinggang Shen, Zhiming Cui: 3D Structure-guided Network for Tooth Alignment in 2D Photograph,15(2), 2023
- 2. Shen Feihong, Liu jingjing, Li Haizhen: ORTHOGAN: HIGH-PRECISION IMAGE GENERATION FOR TEETH ORTHODONTIC VISUALIZATION,15(10),2022
- Lingchen Yang, Zefeng Shi, Yiqian Wu, Xiang Li, Kun Zhou, Hongbo Fu, and Youyi Zheng. iorthopredictor: model-guided deep prediction of teeth alignment. ACM Transactions on Graphics (TOG), 39(6):1–15, 2020.
- Beijia Chen, Hongbo Fu, Kun Zhou, and Youyi Zheng. Orthoaligner: Image-based teeth alignment prediction via latent style manipulation. IEEE Transactions on Visualization and Computer Graphics, 2022.
- Zhiming Cui, Changjian Li, and Wenping Wang. Toothnet: automatic tooth instance segmentation and identification from cone beam ct images. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6368–6377, 2019
- Zhiming Cui, Yu Fang, Lanzhuju Mei, Bojun Zhang, Bo Yu, Jiameng Liu, Caiwen Jiang, Yuhang Sun, Lei Ma, Jiawei Huang, et al. A fully automatic ai system for tooth and alveolar bone segmentation from cone-beam ct images. Nature communications, 13(1):2096, 2022.
- 7. W De Vos, Jan Casselman, and GRJ19464146 Swennen. Cone-beam computerized tomography (cbct) imaging of the oral and maxillofacial region: a systematic review of the literature. International journal of oral and maxillofacial surgery, 38(6):609–625, 2009.
- 8. Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks.

Communications of the ACM, 63(11):139-144, 2020.

- 9. Andreas Wirtz, Florian Jung, Matthias Noll, Angi Wang, and Stefan Wesarg. Automatic modelbased 3-d reconstruction of the teeth from five predefined photographs with viewina directions. In Medical Imaging 2021: Image Processing, volume 11596, pages 198–212. SPIE, 2021.
- 10. Chenglei Wu, Derek Bradley, Pablo Garrido, Michael Zollhöfer, Christian Theobalt, Markus H Gross, and Thabo Beeler. Model-based teeth reconstruction. ACM Trans. Graph., 35(6):220-1, 2016.
- 11. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840-6851, 2020.
- 12. YANG Lingchen, SHI Zefeng, WU Yigian, 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.
- 13. Dor Bank, Noam Koenigstein, Raja Giryes: Autoencoders, 2021.
- 14. Diederik P. Kingma, Max Welling: An Introduction to Variational Autoencoders 2019.
- 15. Guillaume Alain, Yoshua Bengio, Li Yao, Jason Yosinski, Eric Thibodeau-Laufer, Saizheng Zhang, and Pascal Vincent. GSNs: generative stochastic networks. Information and Inference: A Journal of the IMA, 5(2):210-249, 2016.
- 16. Chitwan Saharia, William Chan, Huiwen Chang, Fleet, and Mohammad Norouzi. Palette: Imageto-image diffusion models. In ACM SIGGRAPH 2022 Conference Proceedings, pages 1-10, 2022.
- 17. Hidemichi Kihara, Wataru Hatakeyama, Futoshi Komine, Kyoko Takafuji, Toshiyuki Takahashi, Jun Yokota, Kenta Oriso, and Hisatomo Kondo. Accuracy and practicality of intraoral scanner in prosthodontic research, 64(2):109–113, 2020.
- 18. Vahid Kazemi and Josephine Sullivan. One millisecond face alignment with an ensemble of

conference on computer vision and pattern recognition, pages 1867-1874, 2014.

- 19. Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 815-823, 2015.
- 20. Changgian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In Proceedings of the European conference on computer vision (ECCV), pages 325-341, 2018.
- 21. Richard P McLaughlin and John C Bennett. The dental vto: an analysis of orthodontic tooth movement. Journal of Clinical Orthodontics: JCO, 33(7):394-403, 1999.
- 22. Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In 2016 fourth international conference on 3D vision (3DV), pages 565-571. leee, 2016.
- 23. Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5549-5558, 2020.
- 24. Reuven Rubinstein. The cross-entropy method for combinatorial and continuous optimization. Methodology and computing in applied probability, 1:127-190, 1999.
- Chris Lee, Jonathan Ho, Tim Salimans, David 25. G Power, J Breckon, M Sherriff, and F McDonald. Dolphin imaging software: an analysis of the accuracy of cephalometric digitization and orthognathic prediction. International journal of oral and maxillofacial surgery, 34(6):619-626, 2005.
 - 26. Pascal Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7):1661-1674, 2011
- dentistry: A literature review. Journal of 27. TX Zheng, Shuai Huang, YF Li, and MC Feng. Key techniques for vision based 3d reconstruction: a review. Zidonghua Xuebao, 46:631-652, 2020.
- regression trees. In Proceedings of the IEEE 28. Bo D Tapley and JM Lewallen. Comparison of several numerical optimization methods.

of Journal Optimization Theory and Applications, 1:1–32, 1967.

- 29. Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 1125-1134, 2017.
- 30. Amit Yadav, Anand Gupta, Ms. Aarushi Thusu: Recognition of Sentiment using Deep Neural Network. in International Journal of Trend in Scientific Research and Development (ijtsrd), 2456-6470, 2023.
- 31. Phillip Roe, Kitichai Rungcharassaeng, Joseph Y.K. Kan, Rishi D.Patel, Wayne V. Campagni, James S. Brudvik: The Influence of Upper Lip Length and Lip Mobility on Maxillary Incisal Exposure, 2012
- 32. Mohan Bhuvaneswaran: Principles of smile design, 0.4103/0972-0707.73387,2010
- 33. Karel Zuiderveld. Contrast limited adaptive histogram equalization. Graphics gems, pages 474–485, 1994
- 34. Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 586-595, 2018
- 35. Guodong Wei, Zhiming Cui, Yumeng Liu, Nenglun Chen, Runnan Chen, Guiging Li, and Wenping Wang. Tanet: towards fully automatic tooth arrangement. In Computer Vision-ECCV 2020: 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XV 16, pages 481-497. Springer, 2020.
- 36. Manoj Kumar Sundar and BDS Venkataraman Chelliah. Ten steps to create virtual smile design templates with adobe photoshop cs6. Compendium, 39(3), 2018.
- 37. Vedant Singh, Surgan Jandial, Ayush Chopra, Siddharth Ramesh, Balaji Krishnamurthy, and Vineeth N Balasubramanian. On conditioning the input noise for controlled image generation with diffusion models. arXiv preprint arXiv:2205.03859, 2022.
- 38. Olaf Ronneberger, Philipp Fischer, and Thomas 49. Sachit Menon, Alexandru Damian, Shijia Hu, Brox. U-net: Convolutional networks for

biomedical image segmentation. In Medical Image Computing and ComputerAssisted Intervention-MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18, pages 234-241. Springer, 2015.

- 39. Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In Computer Vision-ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14, pages 694-711. Springer, 2016.
- 40. Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, Sungroh Yoon: ILVR: Conditioning Method for Denoising Diffusion Probabilistic Models, 2021
- 41. Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In ICCV, 2019.
- 42. Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In CVPR, 2020.
- 43. Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. arXiv preprint arXiv:1809.11096, 2018.
- 44. Jooyoung Choi, Jungbeom Lee, Yonghyun Jeong, and Sungroh Yoon. Toward spatially unbiased generative models. arXiv preprint arXiv:2108.01285, 2021.
- 45. Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In ICCV, 2019
- 46. Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In ECCV, 2018.
- 47. Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In CVPR, 2020.
- 48. Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In CVPR, 2019.
 - Nikhil Ravi, and Cynthia Rudin. Pulse: Self-

supervised photo upsampling via latent space exploration of generative models. In CVPR, 2020.

- 50. Taesung Park, Alexei A Efros, Richard Zhang, and JunYan Zhu. Contrastive learning for unpaired image-to-image translation. In ECCV, 2020.
- 51. Olaf Ronneberger, Philipp Fischer, and Thomas Unet: Convolutional networks Brox. for biomedical image segmentation. In International Conference on Medical image computing and computer-assisted intervention, 2015.
- 52. Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In CVPR, 2020.
- 53. Lynton Ardizzone, Carsten Lüth, Jakob Kruse, Carsten Rother, and Ullrich Köthe. 2019. Guided Image Generation with Conditional Invertible Neural Networks. In arXiv:1907.02392.
- 54. Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. 2009. PatchMatch: A Randomized Correspondence Algorithm for Structural Image Editing. ACM Transactions on Graphics (Proc. SIGGRAPH) 28, 2009.
- 55. Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. 2017. Rethinking atrous convolution for semantic image segmentation. arXiv preprint arXiv:1706.05587 (2017).
- 56. Philipp Kopp, Derek Bradley, Thabo Beeler, Landmark Facial Detection, 10.13140/RG.2.2.10980.42886, 2019.
- 57. J. Yang, Q. Liu, and K. Zhang. Stacked Hourglass Network for Robust Facial Landmark Localisation. In 2017 IEEE Conference on Workshops (CVPRW), pages 2025–2033, 2017
- 58. Amir Zadeh, Tadas BaltruÅ ,aaitis, and Louis-Philippe Morency. Convolutional Experts Constrained Local Model for Facial Landmark Detection. arXiv:1611.08657 [cs], November 2016. arXiv: 1611.08657
- 59. Xiangyu Zhu, Zhen Lei, Xiaoming Liu, Hailin Shi, and Stan Z. Li. Face alignment across large poses: A 3d solution. In Proceedings of the IEEE

conference on computer vision and pattern recognition, pages 146-155, 2016

- 60. Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint Face Detection and Alignment using Multi-task Cascaded Convolutional Networks. IEEE Signal Processing Letters, 23(10):1499-1503, October 2016. arXiv: 1604.02878
- 61. Abdul Mueed Hafiz, Ghulam Mohiuddin Bhat: A Survey on Instance Segmentation: State of the art, arXiv:2007.00047, 2020.
- 62. Chien-Yao Wang, Alexey Bochkovskiy, Hong-Yuan Mark Liao: YOLOv7: Trainable bag-offreebies sets new state-of-the-art for real-time object detectors, arXiv:2207.02696, 2022.
- 63. Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. 2018. Learning representations and generative models for 3D point clouds. In International Conference on Machine Learning (ICML), 40-49
- 64. Yazeed Alharbi, Neil Smith, and Peter Wonka. 2019. Latent Filter Scaling for Multimodal Unsupervised Image-to-Image Translation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 1458-1466
- 65. Matthew Amodio and Smita Krishnaswamy. 2019. TraVeLGAN: Image-to-image Translation by Transformation Vector Learning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 8983-8992
- Markus Gross: Analysis and Improvement of 66. Moab Arar, Yiftach Ginger, Dov Danon, Ilya Leizerson, Amit Bermano, and Daniel Cohen-Or. 2020. Unsupervised Multi-Modal Image Registration via Geometry Preserving Image-to-Image Translation. arXiv preprint arXiv:2003.08073,2020.
- Computer Vision and Pattern Recognition 67. Tero Karras, Samuli Laine, Timo Aila: A Style-Based Generator Architecture for Generative Adversarial Networks, arXiv:1812.04948, 2018.
 - 68. Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Muller, Harry Saini Yam Dominik Lorenz, Axel Sauer, Frederic Levi, Boesel, Dustin Podell, Tim Dockhorn Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, Robin Rombach: Scaling Rectified Flow

Transformers for High-Resolution Image Synthesis, arXiv:2403.03206, 2024

- 69. Blattmann, A., Dockhorn, T., Kulal, S., Mendelevitch, D., Kilian, M., Lorenz, D., Levi, Y., English, Z., Voleti, V., Letts, A., et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. arXiv preprint arXiv:2311.15127, 2023
- Brooks, T., Holynski, A., and Efros, A. A. Instructpix2pix: Learning to follow image editing instructions. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 18392–18402, 2023.
- 71. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin: Attentio Is All You Need, arXiv:1706.03762, 2017
- 72. Huiwen Chang Han Zhang Lu Jiang Ce Liu[°] William T. Freeman: MaskGIT: Masked Generative Image Transformer, 2022.