

FACEGPT: SELF-SUPERVISED LEARNING TO CHAT ABOUT 3D HUMAN FACES

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce FaceGPT, a self-supervised learning framework for large vision-language models (VLMs) to reason about 3D human faces from images and text. Typical 3D face analysis algorithms are specialized and lack semantic reasoning capabilities. FaceGPT overcomes this limitation by embedding the parameters of a 3D morphable face model (3DMM) into the token space of a VLM, enabling the generation of 3D faces from both textual and visual inputs. FaceGPT is trained as a model-based autoencoder in a self-supervised manner from in-the-wild images. In particular, a dedicated face token is projected to 3DMM parameters and then rendered as a 2D face image to guide the self-supervised learning process through image-based reconstruction. Without relying on expensive 3D annotations, FaceGPT learns to generate 3D faces based on visual or textual inputs, achieving a competitive performance compared to methods that are specialized to each of these tasks. Most importantly, FaceGPT is able to leverage the world knowledge in VLMs to achieve semantic reasoning capabilities, allowing the model to perform *speculative generation* of 3D faces purely from subtle textual prompts that do not explicitly describe facial features. This opens a new way of generating 3D faces from subtle descriptions of emotions or general everyday situations.

1 INTRODUCTION

In this work, we address the problem of reasoning about 3D human faces from images and text. Related work on monocular 3D face reconstruction aims to estimate the parameters of a 3D morphable model Blanz & Vetter (1999); Tewari et al. (2017); Deng et al. (2019b); Feng et al. (2021a); Li et al. (2023) given 2D face images as input. However, these methods lack the capability to reason about faces from text input. Unlike these systems, humans can vividly imagine and even draw faces based solely on either face images or textual descriptions. Motivated by recent advances in large vision-Language models (VLMs) Liu et al. (2023b); Zhu et al. (2023), we aim to explore a path forward towards enabling VLMs to obtain an in-depth reasoning-based understanding of 3D faces.

To investigate this question, we present FaceGPT, a vision-language model with an intricate ability to reason about 3D human faces from visual and textual input. We represent faces as 3D morphable model (3DMM) parameters that include parameters for the 3D face shape, expression, albedo, and scene illumination. Following related work on image segmentation Lai et al. (2024) and human pose estimation Feng et al. (2024), we extend the original token space of the VLM by incorporating a new `<FACE>` token that is decoded into 3DMM parameters using an MLP (Fig. 1). Thus, when combined with a differentiable computer graphics renderer Ravi et al. (2020), the VLM model becomes capable of synthesizing face images. This enables us to formulate FaceGPT within a model-based autoencoder framework Tewari et al. (2017), and hence to train our model in a fully self-supervised manner from in-the-wild images. To the best of our knowledge, this is the first work combining vision-language model with an inverse graphics pipeline. During training, we freeze the visual encoder of the VLM while training the MLP and the LLM using LoRA Hu et al. (2022). The model is trained with three types of data: (1) In-the-wild face images for the self-supervised training of the `<FACE>` token and 3DMM projection layers via inverse rendering. (2) Text-to-3DMM data for generating 3DMM parameters from text that explicitly describes facial features. (3) Standard multi-modal instruction tuning data to retain the general capability and quality of the VLM. We construct this dataset from a set of face images by running an off-the-shelf self-supervised monocular face reconstruction method Li et al. (2023) and by generating textual descriptions of the faces via the original VLM.

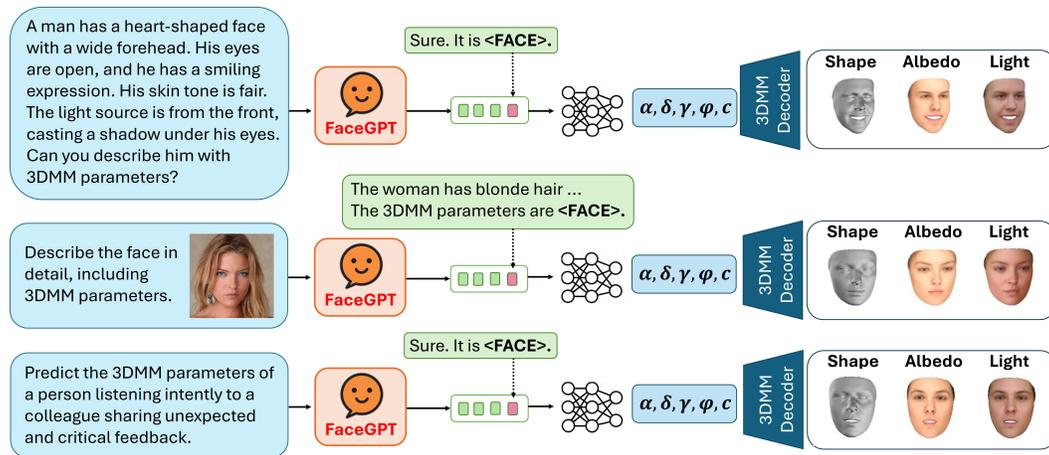


Figure 1: We introduce FaceGPT, a large vision-language model that learns to produce 3D human faces (in terms of 3DMM parameters) in a fully-self-supervised manner. When prompted with face images and task-specific questions, FaceGPT can output a special <FACE> token of which the corresponding feature embedding (red) can be decoded into 3DMM parameters, that encode the face shape α , expression δ , the texture γ , the light ϕ and camera c parameters. When decoded with a 3DMM and differentiable renderer, this enables a fully self-supervised learning via inverse rendering. FaceGPT is a general-purpose model that can produce: 3D human faces from text-only input (first row), as well as from multi-modal input (second row). Moreover, FaceGPT is the first model capable of speculative face generation (last row), all while retaining general chatting abilities.

We evaluate FaceGPT on a variety of tasks, including text-to-3DMM face generation, traditional 3D face reconstruction, and general-purpose visual instruction following. We demonstrate that FaceGPT becomes a general-purpose model that achieves competitive results when compared to specialized methods in all those tasks. Most importantly, we show that FaceGPT is able to leverage the world knowledge in VLMs to achieve semantic reasoning capabilities, allowing the model to perform *speculative generation* of 3D faces purely from subtle textual prompts that do not explicitly describe facial features, such as “the person is listening intently to a colleague sharing unexpected and critical feedback” (Fig. 1). Hence, FaceGPT goes far beyond existing methods as it can translate implicit descriptions of emotional states into 3D facial features, which requires a semantic understanding of (i) how feelings like contemplation affect expressions and (ii) how these changes appear in realistic 3D facial features. Beyond face analysis, we believe that the design principles underlying FaceGPT are general and also suggest a pathway towards a self-supervised integration of the “world knowledge” that VLMs derive from extensive textual data and the structured 3D representations of the visual world via self-supervised learning through inverse rendering. In summary, our work makes the following concrete contributions:

- **A novel vision-language model (FaceGPT) for 3D face reasoning.** We propose the first approach to integrating 3D face understanding capabilities within a vision-language model (VLM). FaceGPT leverages both visual and textual inputs to reason about 3D facial geometry and appearance, bridging the gap between image-based 3D face reconstruction and text-based facial description interpretation.
- **Semantic reasoning for speculative 3D face generation.** Unlike traditional methods that rely on explicit visual or textual cues, FaceGPT demonstrates the ability to perform speculative 3D face generation based on abstract or emotional descriptions.
- **Competitive performance across multiple tasks.** We evaluate FaceGPT across various benchmarks, including 3D face reconstruction from images, text-to-3D face generation, and visual instruction following. Our results show that FaceGPT performs competitively with specialized 3D face reconstruction methods while maintaining the flexibility and reasoning capabilities of a general-purpose vision-language model.

108 The design principles of FaceGPT, specifically the integration of structured 3D representations
109 with the world knowledge encoded in VLMs, provide a general framework for learning 3D-aware
110 multimodal reasoning from in-the-wild 2D images in a self-supervised manner. Therefore, we believe
111 FaceGPT represents a significant step forward in the field of vision-language models.

113 2 RELATED WORKS

115 2.1 MONOCULAR MODEL-BASED FACE RECONSTRUCTION

117 Realistically reconstructing digital human faces has been a longstanding challenge in computer vision
118 and graphics due to their vast potential applications. Traditional methods primarily use parametric
119 3D Morphable Models (3DMM) Blanz & Vetter (1999); Paysan et al. (2009); Li et al. (2017) with
120 PCA for dimensionality reduction to simplify high-dimensional 3D face scans, serving as a 3D
121 prior for representing unique facial characteristics and providing precise control. Recently, deep
122 learning-based methods that map 2D images to 3D face models have gained popularity. Early
123 methods struggled with the need for extensive 3D facial scan data paired with 2D images, which
124 was labor-intensive and costly. This limitation was addressed with the introduction of model-based
125 face autoencoders (MoFA) Tewari et al. (2017) that enabled self-supervised 3D face reconstruction.
126 MoFA uses a differentiable rendering layer to minimize differences between input and rendered
127 images, enabling end-to-end learning without ground-truth 3D faces, leading to a number of effective
128 extensions in the self-supervised learning strategies Tewari et al. (2018b;a); Bas et al. (2017); Genova
129 et al. (2018); Daněček et al. (2022). RingNet Sanyal et al. (2019) and DECA Feng et al. (2021b)
130 use landmark-based training, predicting landmarks for input images and treating them as pseudo
131 ground truth, measuring the distance between 2D face landmarks and their projections on the 3DMM
132 surface. The FOCUS Li et al. (2023) framework jointly trains a face autoencoder and an outlier
133 segmentation network, which makes the method robust to outliers such as occlusion and make-up.
134 These advancements significantly improved monocular model-based face reconstruction, making
135 it more efficient and effective. However, these methods are highly specialized and lack a deep
136 understanding of the semantics of human faces or the ability to relate faces to language, limiting their
137 overall scope and effectiveness.

138 2.2 TEXT-TO-3D FACE GENERATION AND MANIPULATION

139 Text-to-3D face generation and manipulation methods aim to use textual information for creating and
140 editing 3D faces. Methods like Dreamface Zhang et al. (2023) and Describe3D Wu et al. (2023b)
141 generate text-conditioned texture maps to render 3D morphable models (3DMM). TG-3DFace Yu
142 et al. (2023) advances this by using tri-plane neural representations and extending the 3D-aware
143 GAN, EG3D Chan et al. (2022), for end-to-end text-conditioned generation. For text-guided 3D
144 face manipulation, methods like Latent3D Canfes et al. (2022) and ClipFace Aneja et al. (2023)
145 optimize intermediate layers with a CLIP-based loss to generate UV-texture maps or predict texture
146 and expression latent codes. These methods, however, rely on 3D scan data and to train new mappers
147 for each text instruction. Unlike these task-specific approaches, FaceGPT is a general-purpose model
148 that reasons about 3D human faces from images, text, or both by leveraging general visual knowledge.
149 Our model can interact with users through conversations, discussing facial features and providing
150 relevant responses, while also being capable of following general user instructions.

152 2.3 MULTIMODAL LARGE LANGUAGE MODELS

153 Large Language Models (LLMs) are rapidly transforming various fields Radford et al. (2019);
154 Brown et al. (2020); OpenAI (2024); OpenAI et al. (2024). While proprietary models like OpenAI’s
155 ChatGPT OpenAI (2024) and GPT-4 OpenAI et al. (2024) dominate the landscape, open-source
156 alternatives such as Vicuna Chiang et al. (2023), LLaMA Touvron et al. (2023), Alpaca Taori et al.
157 (2023), Mistral Jiang et al. (2023) and Qwen Bai et al. (2023) support research efforts. However,
158 LLMs mainly focus on generating text as output given text-only input. The integration of additional
159 modalities into LLMs represents an active area of research.

160 Multi-Modal Large Language Models (MM-LLMs) are emerging, extending LLMs’ capabilities be-
161 yond text to encompass a broader spectrum of modalities, including images, videos, and audio. In the

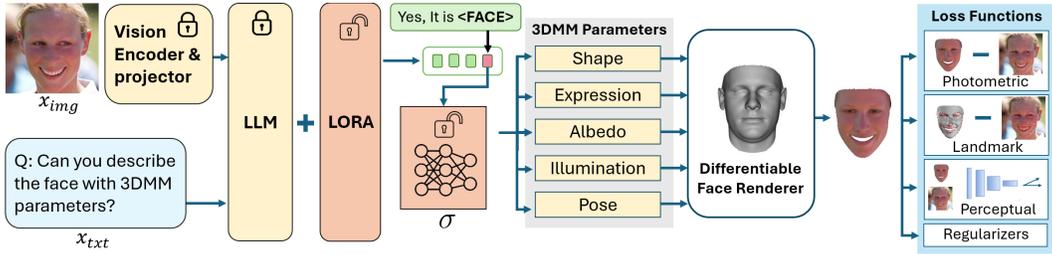


Figure 2: Architecture of FaceGPT. Our model consists of a vision-language model, which includes a vision encoder, a vision projection layer, and an LLM, along with a 3DMM projection layer, denoted as σ , and the parametric Basel face model Blanz & Vetter (1999). During training, the σ projection layer is optimized and the LLM is fine-tuned through LORA, while keeping other components frozen. The training is guided through a self-supervised reconstruction loss using a differentiable renderer.

realm of image-text understanding, recent endeavors like LLaVA Liu et al. (2024) and MiniGPT-4 Zhu et al. (2023) incorporate vision encoders to interpret images and align their features with language embeddings using projection layers. Moreover, cutting-edge models such as PandaGPT Su et al. (2023), ImageBind Girdhar et al. (2023), and NeXT-GPT Wu et al. (2023c) exhibit impressive versatility in handling diverse modalities, aligning embeddings from text, images, audio, and video with language as both input and output. To enhance LLM to natively output more modalities, approaches like LISA Lai et al. (2024) connects LLaVA with a decoder to generate text and segmentation masks, while ChatPose Feng et al. (2024) specializes in human pose information. However, these methods typically rely on supervised learning.

Our aim is to develop a general-purpose vision-language model that is able to (1) generate 3D faces from text or image inputs and (2) capable of connecting the world knowledge from the LLM with 3D human faces to achieve semantic reasoning capabilities about faces in a self-supervised manner.

3 METHOD

In this work, our aim is to augment existing large Vision Language Models (VLMs) with the ability of reasoning about 3D human faces without requiring manual human efforts. Inspired by established model-based face reconstruction methods, we represent 3D human face with 3D Morphable Model (3DMM) Blanz & Vetter (1999) parameters, representing the *face shape*, *expression*, *albedo*, *illumination* and *pose*. In particular, we introduce a $\langle \text{FACE} \rangle$ token into the language space of the LLM, which is mapped into the 3DMM parameter space and subsequently rendered into a 2D image, hence enabling self-supervised 3D facial reconstruction. Fig. 2 presents the whole pipeline of FaceGPT.

3.1 MODEL ARCHITECTURE

Representing Face in language space. Inspired from recent advancements in LMMs, we treat the human face as a distinct modality and incorporate its representation into the language space of VLM. Specifically, we extend the vocabulary of VLM to include a new token $\langle \text{FACE} \rangle$ that specifically represents the human face. Given an input text prompt \mathbf{x}_{txt} and/or input image \mathbf{x}_{img} , the VLM f predicts text responses:

$$\mathbf{y}_{txt} = f(\mathbf{x}_{img}, \mathbf{x}_{txt}), \quad (1)$$

where $\mathbf{y}_{txt} = [t_1, \dots, t_N]$ is the output sequence of tokens with corresponding hidden states $[h_1, \dots, h_N]$. When \mathbf{x}_{txt} contains a face generation instruction, the resulting output \mathbf{y}_{txt} will include a $\langle \text{FACE} \rangle$ token, facilitating further 3DMM predictions.

From $\langle \text{FACE} \rangle$ token to 3DMM. If one of the output tokens $t_n \in \mathbf{y}_{txt}$ is our defined $\langle \text{FACE} \rangle$ token, we can extract the hidden state as $h_{\langle \text{FACE} \rangle} = h_n \in \mathbb{R}^{4096}$ and project it using an MLP σ into the latent 3DMM parameters $\theta = \sigma(h_{\langle \text{FACE} \rangle}) = [\alpha, \delta, \gamma, \phi, c] \in \mathbb{R}^{257}$, i.e. the 3D face shape $\alpha \in \mathbb{R}^{80}$, facial expression parameters $\delta \in \mathbb{R}^{64}$ and texture $\gamma \in \mathbb{R}^{80}$ of a 3DMM, as well as the spherical harmonics illumination $\phi \in \mathbb{R}^{27}$ and camera parameters $c \in \mathbb{R}^6$ of the scene. The 3D vertices and triangles, as well as the color of the face mesh are then determined using the standard 3DMM model

$M(\theta)$ as described in Tewari et al. (2017). Using an orthographic camera model the reconstructed 3D face mesh can be rendered into 2D space using a differentiable renderer Π , hence producing the final reconstructed face image \hat{y}_{rec} . The process can be summarized as:

$$\theta = \sigma(h_{\langle \text{FACE} \rangle}) \quad (2a)$$

$$\hat{y}_{rec} = \Pi(M(\theta)) \quad (2b)$$

Note that the FaceGPT architecture can be seen as a new type of language-based autoencoder, with a VLM as encoder and a computer graphics decoder that is based on the 3DMM.

3.2 SELF-SUPERVISED TRAINING

Our objective is to develop a VLM that is able to learn an enhanced comprehension of 3D human faces in a self-supervised manner. Ultimately, the model should not only understand user instructions reliably, but also to accurately reconstruct 3D faces from visual or text input. To accomplish this, we have devised a self-supervised approach that incorporates 3DMM understanding into existing VLMs. This method allows us to leverage face data without the need for costly 3DMM annotations. We also construct a text-to-3DMM dataset in an unsupervised manner that supports this training paradigm, enabling our model to learn these capabilities effectively and efficiently.

Self-supervised face reconstruction loss. To incorporate the 3DMM as a new modality into an existing VLM, we add a new token $\langle \text{FACE} \rangle$ in the vocabulary of the LLM and fine-tune the language modelling head. For fine-tuning the VLM without relying on manually annotated data, we introduce a 2D self-supervised loss L_{face} , following established protocols of specialized face reconstruction models Li et al. (2023):

$$L_{face} = \lambda_{pixel}L_{pixel} + \lambda_{perc}L_{perc} + \lambda_{LM}L_{LM} + \lambda_{reg}L_{reg}, \quad (3)$$

Specifically, L_{face} incorporates the following components:

Reconstruction loss. $L_{pixel} = \|A \odot (x_{img} - \hat{y}_{rec})\|_2^2$ refers to the pixel-wise reconstruction loss between reconstructed images \hat{y}_{rec} and the input images x_{img} . To avoid distortions when training from in-the-wild data, a 2D skin mask A is estimated from the input images using a simple pre-trained Gaussian mixture model for the skin color Deng et al. (2019b).

Perceptual loss. $L_{perc} = sim(f_{perc}(x_{img}), f_{perc}(\hat{y}_{rec}))$ estimates the cosine similarity between the reconstructed image \hat{y}_{rec} and the input image x_{img} at the perceptual level using a pre-trained feature extractor f_{perc} .

Landmark loss. $L_{LM} = \|LM_{img} - LM_{rec}\|_2^2$ measures the L2 distance between the projected 2D landmarks of the estimated 3DMM LM_{rec} and predicted 2D facial landmarks LM_{img} using an off-the-shelf facial detector (Bulat & Tzimiropoulos, 2017).

Parameter regularization. $L_{reg} = \|\theta\|_2^2$ regularizes the estimated 3DMM parameters towards the mean value of the multi-variate Gaussian distribution of the 3DMM.

The weight parameters $\lambda_{face} = [\lambda_{pixel}, \lambda_{perc}, \lambda_{LM}, \lambda_{reg}]$ balance the respective losses, and we set them as established in prior work Li et al. (2023).

Preserving the ability for natural conversations about faces. We noticed that VLM’s lose the ability to conduct general conversations about faces when trained self-supervised face reconstruction loss, and tend to always output 3DMM parameters when queried with a face image. To resolve this problem, we generate a face conversation dataset with accurate textual face descriptions, by querying the VLM with face images and conversations that include multiple questions about facial attributes. During training, we mix task-specific instructions that explicitly ask for 3DMM parameters with general conversational data to regularize the training and preserve the ability for general non-3DMM related conversations about faces. In contrast to prior works (Lai et al., 2024; Feng et al., 2024) which generally employ a simple static question-answer template that restrains the model from providing diverse responses, our proposed strategy can effectively improve the instruction following performance for general face conversations.

Self-supervised Text-to-3DMM loss. We aim to enable the VLM to perform semantic face understanding tasks, i.e. to predict faithful 3DMM parameters when only having textual input $f(\cdot, x_{txt})$.

While related works for 3D human pose estimation Delmas, Ginger and Weinzaepfel, Philippe and Lucas, Thomas and Moreno-Noguer, Francesc and Rogez, Grégory (2022); Feng et al. (2024) rely on human-written fine-grained language descriptions of 3D human meshes, we aim to achieve text-to-3DMM generation without any human interaction. To achieve this, we introduce a self-supervised text-based face reconstruction loss. In particular, we first query the VLM with the face images in the training data and a question template that requests detailed face descriptions of a number of facial attributes such as head shape, skin tone, or facial expressions. This results in a dataset with pairs of face images and detailed textual descriptions about facial attributes. During training, we instruct the VLM to estimate 3DMM parameters only from the available detailed face descriptions $f(\cdot, x_{txt})$. As we have the corresponding 2D image for each face description, we can re-use the self-supervised face reconstruction loss Eq. (3) to train the model at the task of text-to-3DMM generation, avoiding the need for any human-written annotation.

Text prediction loss. We also utilize the autoregressive objective L_{txt} to guide the model to produce correct text output, thereby preserving its general ability to follow user instructions:

$$L_{txt} = CE(\hat{y}_{txt}, y_{txt}), \quad (4)$$

where y_{txt} is the ground truth text response and CE is the cross entropy loss.

Overall training Objective. Our overall training objective L integrates both text-based autoregressive objective and a self-supervised loss for predicting 3DMM from either image or text input:

$$L = \lambda_{txt}L_{txt} + \lambda_{face}L_{face} \quad (5)$$

where λ_{txt} and λ_{face} are the weights for balancing the losses.

3.3 SEMANTIC REASONING ABOUT HUMAN FACES

After being trained as described in the previous section, FaceGPT becomes a general model that can estimate 3D facial expressions from single images, create facial features based on detailed descriptions, and participate in question-and-answer dialogues. Most remarkably, even without directly training the model for linking 3D facial attributes to subtle phrases that do not directly contain specific facial traits, our model shows a zero-shot ability to reason about human faces from descriptions that contain emotions or general descriptions of everyday situations. This means the model can combine reasoning and world knowledge with the 3D facial representation. Similar as in ChatPose Feng et al. (2024), our aim is to highlight these emerging capabilities by introducing the task of *speculative face generation*, which focuses on the model’s ability to reason about 3D human faces from speculative queries.

Speculative Face Generation. In this task, rather than providing explicit facial descriptions that directly detail the shape and texture of features, we pose speculative queries related to a person’s emotional state or general everyday situations. The model is then tasked to infer a plausible 3D face, based on the assumption that the person is experiencing the described situation or emotion. For instance, a user might say, “Predict the face of a person who is excited about a surprise party.” Answering such queries requires an understanding of broader concepts like “excitement” and the ability to deduce the appropriate facial features, followed by generating the relevant 3D facial parameters. To build an evaluation dataset, we draw from the CelebA Text dataset Sun et al. (2021) for facial descriptions. We then use GPT-4 to rephrase these descriptions into questions about the emotions tied to each expression, resulting in a total of 444 responses, with 64 examples selected for evaluation. These responses undergo manual review and corrections in order to remove any direct descriptions of facial features. We further estimate the 3DMM face parameters of every test image using a state-of-the-art face autoencoder Li et al. (2023), resulting in paired data with speculative descriptions and corresponding 3D face parameters. More details can be found in supplementary.

4 EXPERIMENTS

4.1 TRAINING DATA

Text-to-Face Data. The text-to-face dataset comprises pairs of text descriptions and face images, facilitating the development of mappings by VLM between textual descriptions and

324 3DMM parameters of faces. During training, only text is taken as input and the correspond-
 325 ing face images are only used for loss computation. Given the absence of publicly available
 326 datasets linking text descriptions to 3D face meshes and existing VLMs like LLaVA present
 327 powerful Visual Question Answering(VQA) ability, we rely on pre-trained VLMs to generate
 328 textual descriptions for face images. We employ following templates to guide the learning
 329 process in VLMs: "USER: {description}, can you give the 3DMM parameters
 330 of this person. ASSISTANT: Sure, it is <FACE>.", where {description}
 331 is the text description for faces. We utilize high-quality face images from CelebA-HQ (Karras et al.,
 332 2018) and use LLaVA (Liu et al., 2023a) to produce explicit text descriptions for dataset construction.

333 **Image-to-Face Data.** Image-to-face reconstruction data is composed of only human face images. The
 334 face images will be formatted with a template like "USER: <IMAGE> Can you give the
 335 3DMM parameters of this person. ASSISTANT: Sure, it is <FACE>".
 336 To enhance the diversity and relevance of the conversations centered around human face images,
 337 we also generate face-centric conversations for each face image. This approach enriches the
 338 contextual understanding of the VLM concerning the newly introduced token <FACE>. We adopt
 339 the CelebA-HQ trainset as the image to face reconstruction dataset.

340 **Multimodal Instruction-Following Data.** This data is general-purpose VQA data, and it is used to
 341 preserve the VLM’s ability of understanding a user’s instructions. Following LLaVA v1.5, We use
 342 LLaVA-v1.5-mix665k as multimodal instruction following data.

344 4.2 EXPERIMENTAL SETTINGS

345 **Network Architecture.** We build our model on Large Vision Language Model LLaVA-1.5-7B (Liu
 346 et al., 2023a) with CLIP-ViT-L-336px as vision encoder and Vicuna v1.5 as the LLM backbone.
 347 LoRA (Hu et al., 2022) is applied to efficiently fine-tune the VLM. An MLP head with GeLU
 348 activations (Hendrycks & Gimpel, 2016) and channels [5120, 5120, 257] is appended to the last layer
 349 of the VLM to predict the 3D human face parameters.

350 **3D Face Model.** We use the Basel Face Model (BFM) 2017 Gerig et al. (2018) as the 3D face model.
 351 The face is parameterized as the semantic code vector $\theta = [\alpha, \delta, \gamma, \phi, c] \in \mathbb{R}^{257}$ in (Tewari et al.,
 352 2017), which includes 3D shape parameters $\alpha \in \mathbb{R}^{80}$, facial expression parameters $\delta \in \mathbb{R}^{64}$, texture
 353 of 3DMM $\gamma \in \mathbb{R}^{80}$, illumination $\phi \in \mathbb{R}^{27}$ and camera parameters $c \in \mathbb{R}^6$.

354 **LLaVa-Key baseline.** As there are no VLM-based methods that can perform 3D face reconstruction,
 355 we introduce the LLaVa-Key. The model is finetuned to predict facial landmarks in the format of pure
 356 text given either images or textual description of human faces as input. For each input, gradient-based
 357 optimization is applied to these predicted landmarks to fit the 3DMM using Eq. (3) to obtain a 3D face
 358 reconstruction. It is important to note that LLaVa-Key employs test-time fine-tuning, making it an
 359 inherently different baseline compared to all other feed-forward methods. Nevertheless, LLaVa-Key
 360 is useful to show the effect of a directly using the token space of VLMs to encode 3D information.

361 **Implementation Details.** The training uses 8 NVIDIA 48G A40 GPUs. We utilize deepspeed (Rasley
 362 et al., 2020) engine and ZeRO optimizer (Rajbhandari et al., 2020) for efficient training. We use
 363 AdamW (Loshchilov & Hutter, 2019) optimizer with learning rate and weight decay set to $2e - 5$ and
 364 0, respectively. We also follow the standard setting of using a WarmupDecayLR as the learning rate
 365 scheduler, where the warm-up iterations are set to 100. The weights of the text generation loss λ_{txt}
 366 and the face reconstruction loss λ_{face} are set to 1.0 and 0.1, respectively. Following (Li et al., 2023),
 367 those of the pixel loss λ_{pixel} , the perceptual loss λ_{per} , the landmark loss λ_{LM} , and the regularization
 368 loss λ_{reg} are set to 0.5, 0.25, $5e-4$ and 0.1, respectively. The landmarks of human faces are obtained
 369 with (Bulat & Tzimiropoulos, 2017) and pre-trained ArcFace (Deng et al., 2019a) is used to compute
 370 perceptual loss. The batch size is set to 8 per GPU and gradient accumulation step is set to 4.

371 **Evaluation Metrics.** For text-based face generation task, we compare the estimated 3D human face
 372 with the ground truth mesh to measure the Chamfer Distance(CD), Complete Rate(CR) and Relative
 373 Face Recognition Rate(RFRR) following Describe3D (Wu et al., 2023a). Both Chamfer Distance and
 374 Complete Rate would be used to reflect the accuracy of 3D meshes. And Relative Face Recognition
 375 Rate can measure the identity similarity of the textured 3D face rendering. We also additionally
 376 perform a user study for speculative face generation task as it is a difficult task requiring common
 377 sense of human. For image-based face reconstruction task, we measure the L2 photometric error
 in RGB space and the L2 landmark error to reflect the quality of 3D face’s eometry and texture.

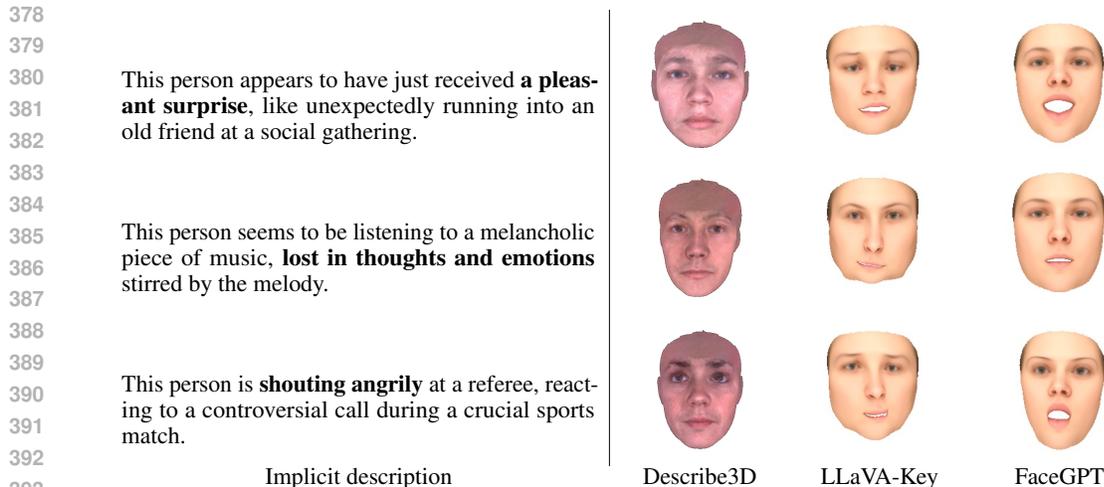


Figure 3: Qualitative results for speculative face generation. The abstract concepts and human activities in the descriptions are highlighted with bold text. FaceGPT presents a better capability in understanding the abstract concept and human activities compared to other methods.

Table 1: Quantitative results for speculative face generation. FaceGPT achieves a significant advantage in terms of the predicted 3D mesh’s accuracy and user preference compared to other methods.

| Method | Unsupervised | CD ↓ | CR (%) ↑ | RFRR ↑ | User Study (%) ↑ |
|------------|--------------|-------------|-------------|-------------|------------------|
| Describe3D | ✗ | 153.1 | 25.6 | 14.0 | 11.2 |
| LLaVA-Key | ✓ | 38.5 | 68.1 | 40.0 | 31.5 |
| FaceGPT | ✓ | 11.5 | 83.6 | 64.0 | 57.2 |

The instruction following ability is measured with GPT-assisted evaluation as described in Liu et al. (2023b), which queries GPT4 to obtain the grading of generated responses.

User Study. For evaluating the quality of speculative face generation task, we also perform a user study with 23 volunteers. Each volunteer will be presented with 20 questions from SPG benchmark and each question consists of an implicit description with the visual results generated by different methods given this description as input. For each question, the volunteer will be asked to utilize their understanding about the abstract descriptions and to select the result which best matches the implicit descriptions. The ratio of the number of times a specific method is selected to the total number of questions will be reported as the result.

4.3 SPECULATIVE FACE GENERATION

In this section, we evaluate FaceGPT’s zero-shot capabilities at speculative face generation. We use the same template in Text-to-Face Data to query model and replace the `{description}` with the implicit description in the benchmark. For a fair comparison, we only select the face region of each method during the evaluation and align the 3D point clouds with the Iterative Closest Point (ICP) method before evaluation. The results are presented in Table 1 quantitatively and in Figure 3 qualitatively. We can observe that FaceGPT has significant advantages compared to Describe3D and the baseline LLaVA-Key method in terms of the quality of 3D shape and in terms of preference ratings in the user study. Interestingly, our baseline LLaVA-Key also outperforms Describe3D which requires a CLIP model and 3D supervision during training. In summary, FaceGPT develops a common sense about human faces that enables it to infer facial features from abstract and indirect descriptions.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

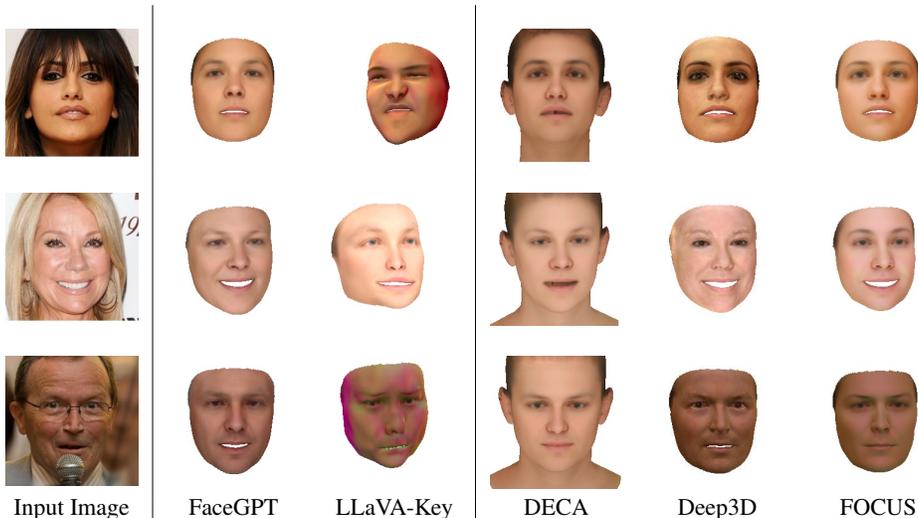


Figure 4: Qualitative results for 3D face reconstruction. Our approach allows for the regression of pose, shape, expression, skin reflectance, and illumination from a single monocular image, achieving quality comparable to recent state-of-the-art methods.

4.4 EXPLICIT TEXT-BASED 3D FACE GENERATION

Table 2 shows the results of FaceGPT at 3D face reconstruction with explicit text description. As there are no public unsupervised methods available that perform this task, we compare FaceGPT to a supervised method Describe3D and our LLaVA-Key baseline. LLaVA-Key would fine-tune LLaVA to predict the 2D coordinates of face landmarks given text description as input and an optimization-based method is utilized to fit a 3DMM on predicted landmarks.

Like in speculative face generation benchmark, We also observe the significant advantage that FaceGPT achieves over other methods. Despite being trained in a self-supervised manner, FaceGPT can achieve faithful text-based 3D face reconstructions. These results demonstrate embedding 3D knowledge about our world into a VLM is in principle possible without detailed human annotations, hence demonstrating the large potential of combining VLM’s with vision-as-inverse graphics for self-supervised learning. The qualitative results can be found in the supplementary.

4.5 IMAGE-BASED 3D FACE RECONSTRUCTION

We evaluate FaceGPT on image-based 3D face reconstruction and compare it with SOTA specialized methods for unsupervised monocular face reconstruction and a VLM-based baseline LLaVa-Key in Table 3.

Moreover, Fig. 4 shows a qualitative comparison of FaceGPT and all baselines at the classic task of 3D face reconstruction from a single image. As reported in prior works for VLM-based segmentation Lai et al. (2024) and human pose estimation Feng et al. (2024), our FaceGPT model does, as expected, not reach the state-of-the-art performance of specialized models. However,

Table 2: Performance of explicit text-based 3D face reconstruction. The evaluation is performed on shape and expression by comparing pseudo ground truth produced by FOCUS method.

| Method | CD ↓ | CR ↑ | RFRR ↑ |
|------------|-------|------|--------|
| Describe3D | 96.88 | 29.3 | 16.27 |
| LLaVA-Key | 16.89 | 71.8 | 42.38 |
| FaceGPT | 7.28 | 91.7 | 64.88 |

Table 3: Performance at classical monocular 3D face reconstruction.

| Method | photo ↓ | keypoint ↓ |
|-----------|---------|------------|
| DECA | 0.216 | 5.2px |
| Deep3D | 0.073 | 3.2px |
| FOCUS | 0.077 | 2.2px |
| LLaVa-Key | 0.110 | 14.0px |
| FaceGPT | 0.103 | 3.0px |

Table 4: GPT4-Assisted Evaluation on instruction-following capability. “Conv”, “Details”, and “Complex” correspond to three types of questions (conversation, detailed description, complex reasoning) produced by LLaVA’s data generation pipeline. GPT4 will be prompted to evaluate the answers from different models along with the ground truth answer produced by text-only GPT4 (gpt-4-0613). It would then give a score for each answer with an explanation.

| | Conv | Detail | Complex | All |
|------------------------------------|------|--------|---------|------|
| ChatPose (Feng et al., 2024) | 74.5 | 81.0 | 93.3 | 82.9 |
| LLaVA-V1.5-13B (Liu et al., 2023a) | 80.4 | 81.4 | 90.9 | 84.2 |
| LLaVA-V1.5-7B (Liu et al., 2023b) | 79.9 | 77.6 | 92.4 | 83.4 |
| FaceGPT | 79.6 | 81.5 | 92.6 | 84.6 |

Table 5: Influence of face-centric conversation generation

| Method | face convs | photo | keypoint | Conv | Detail | Complex | All |
|--------------|------------|-------|----------|------|--------|---------|------|
| LLaVA-1.5-7B | ✗ | - | - | 79.9 | 77.6 | 92.4 | 83.4 |
| FaceGPT | ✗ | 0.110 | 3.2px | 78.4 | 80.8 | 89.0 | 82.7 |
| FaceGPT | ✓ | 0.103 | 3.0px | 79.6 | 81.5 | 92.6 | 84.6 |

we note that FaceGPT has a frozen vision encoder to preserve the generalist behavior, whereas all baseline models have a fine-tuned task-specific backbone. Moreover, our model does outperform DECA, which is a highly competitive and widely applied baseline model Zheng et al. (2023). When compared to the VLM baseline model, FaceGPT achieves large improvements, highlighting the benefit of embedding the 3DMM parameters directly in the token space of the VLM.

4.6 GPT-ASSISTED EVALUATION

Table 4 shows that FaceGPT preserves the ability of instruction following by following LLaVA’s evaluation (Liu et al., 2023b) protocol, using GPT4-Assisted evaluation on LLaVA-Bench (COCO). FaceGPT compares favorably to LLaVA-v1.5-7B and reaches similar or better performance in the benchmark compared to VLMs with more parameters. This performance advantage can be attributed to the specialized face training dataset and the usage of our face-centric conversation. These elements enhance FaceGPT’s proficiency in interpreting face-related language-based instructions, improving its overall effectiveness in relevant tasks.

4.7 ABLATION STUDY

Influence of face-centric conversation generation. To prove the necessity of our face-centric conversation generation strategy, we train a model only using a simple single-turn conversation template for face images, which is a common strategy used in the VLM-based image understanding works like LISA and ChatPose. The comparison results are demonstrated in Table 5. We observe that face-centric conversations help a lot in improving model’s ability in performing detailed description and complex reasoning. The face-centric conversations only have a small effect on the face reconstruction task, which is expected, as the generated conversations do not contain information about the 3DMM parameters of faces.

5 CONCLUSION

FaceGPT is the first self-supervised learning framework for Large Vision-Language Models to reason about 3D human faces. We show that VLMs can learn to predict detailed 3D human faces from not only images, but also from textual inputs, in a fully self-supervised manner via inverse rendering. As a generalist model, FaceGPT achieves strong results across various tasks, including text-based face generation, traditional 3D face reconstruction, visual instruction following. FaceGPT also presents impressive ability to infer 3D faces from abstract and indirect text descriptions. We believe our work also has general implications beyond face analytics, as it points towards a way forward to enable large multi-modal language-models to reason about our 3D world without supervision.

REFERENCES

- 540
541
542 Shivangi Aneja, Justus Thies, Angela Dai, and Matthias Nießner. Clipface: Text-guided editing of
543 textured 3d morphable models. In *ACM SIGGRAPH 2023 Conference Proceedings*, pp. 1–11,
544 2023.
- 545 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
546 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu,
547 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan,
548 Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin
549 Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng
550 Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou,
551 Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*,
552 2023.
- 553 Anil Bas, Patrik Huber, William AP Smith, Muhammad Awais, and Josef Kittler. 3d morphable
554 models as spatial transformer networks. In *Proceedings of the IEEE International Conference on*
555 *Computer Vision Workshops*, pp. 904–912, 2017.
- 557 V Blanz and T Vetter. A morphable model for the synthesis of 3d faces. In *26th Annual Conference*
558 *on Computer Graphics and Interactive Techniques (SIGGRAPH 1999)*, pp. 187–194. ACM Press,
559 1999.
- 560 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-
561 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
562 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,
563 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-
564 teusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-
565 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot
566 learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Ad-
567 vances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Asso-
568 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/
569 2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- 570 Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment
571 problem? (and a dataset of 230,000 3d facial landmarks). In *International Conference on Computer*
572 *Vision*, 2017.
- 574 Zehranaz Canfes, M. Furkan Atasoy, Alara Dirik, and Pinar Yanardag. Text and image guided 3d
575 avatar generation and manipulation, 2022.
- 576 Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio
577 Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d
578 generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision*
579 *and pattern recognition*, pp. 16123–16133, 2022.
- 581 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
582 Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot
583 impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April
584 2023), 2(3):6, 2023.
- 585 Radek Daněček, Michael J Black, and Timo Bolkart. Emoca: Emotion driven monocular face capture
586 and animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
587 *Recognition*, pp. 20311–20322, 2022.
- 588 Delmas, Ginger and Weinzaepfel, Philippe and Lucas, Thomas and Moreno-Noguer, Francesc and
589 Rogez, Grégory. PoseScript: 3D Human Poses from Natural Language. In *ECCV*, 2022.
- 591 Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin
592 loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
593 *and Pattern Recognition*, pp. 4690–4699, 2019a.

- 594 Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face
595 reconstruction with weakly-supervised learning: From single image to image set. In *IEEE*
596 *Computer Vision and Pattern Recognition Workshops*, 2019b.
- 597 Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. Learning an animatable detailed 3D
598 face model from in-the-wild images. *ACM Transactions on Graphics (ToG), Proc. SIGGRAPH*, 40
599 (4):88:1–88:13, August 2021a.
- 600 Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart. Learning an animatable detailed 3d
601 face model from in-the-wild images. *ACM Transactions on Graphics (ToG)*, 40(4):1–13, 2021b.
- 602 Yao Feng, Jing Lin, Sai Kumar Dwivedi, Yu Sun, Priyanka Patel, and Michael J. Black. Chatpose:
603 Chatting about 3d human pose. In *CVPR*, 2024.
- 604 Kyle Genova, Forrester Cole, Aaron Maschinot, Aaron Sarna, Daniel Vlasic, and William T Freeman.
605 Unsupervised training for 3d morphable model regression. In *Proceedings of the IEEE conference*
606 *on computer vision and pattern recognition*, pp. 8377–8386, 2018.
- 607 Thomas Gerig, Andreas Morel-Forster, Clemens Blumer, Bernhard Egger, Marcel Luthi, Sandro
608 Schönborn, and Thomas Vetter. Morphable face models-an open framework. In *2018 13th IEEE*
609 *international conference on automatic face & gesture recognition (FG 2018)*, pp. 75–82. IEEE,
610 2018.
- 611 Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand
612 Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the*
613 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15180–15190, 2023.
- 614 Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian
615 error linear units. *CoRR*, abs/1606.08415, 2016. URL [http://arxiv.org/abs/1606.](http://arxiv.org/abs/1606.08415)
616 [08415](http://arxiv.org/abs/1606.08415).
- 617 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
618 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International*
619 *Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=nZeVKeeFYf9)
620 [id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).
- 621 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
622 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
623 L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas
624 Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023. URL [https://arxiv.](https://arxiv.org/abs/2310.06825)
625 [org/abs/2310.06825](https://arxiv.org/abs/2310.06825).
- 626 Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for
627 improved quality, stability, and variation. 2018.
- 628 Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning
629 segmentation via large language model. In *CVPR*, 2024.
- 630 Chunlu Li, Andreas Morel-Forster, Thomas Vetter, Bernhard Egger, and Adam Kortylewski. Robust
631 model-based face reconstruction through weakly-supervised outlier segmentation. In *Proceedings*
632 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 372–381, 2023.
- 633 Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial
634 shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017.
- 635 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
636 tuning, 2023a.
- 637 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*,
638 2023b.
- 639 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in*
640 *neural information processing systems*, 36, 2024.

- 648 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*
649 *ence on Learning Representations*, 2019. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Bkg6RiCqY7)
650 [Bkg6RiCqY7](https://openreview.net/forum?id=Bkg6RiCqY7).
651
- 652 OpenAI. ChatGPT: A conversational language model. Online, 2024. URL [https://www.](https://www.openai.com/chatgpt)
653 [openai.com/chatgpt](https://www.openai.com/chatgpt).
- 654 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
655 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor
656 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian,
657 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny
658 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks,
659 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea
660 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,
661 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung,
662 Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,
663 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty
664 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte,
665 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel
666 Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua
667 Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike
668 Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon
669 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne
670 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo
671 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar,
672 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik
673 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich,
674 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy
675 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie
676 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini,
677 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne,
678 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David
679 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie
680 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély,
681 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo
682 Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano,
683 Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng,
684 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto,
685 Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power,
686 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis
687 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted
688 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel
689 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon
690 Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
691 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
692 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston
693 Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,
694 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason
695 Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff,
696 Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu,
697 Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba,
698 Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang,
699 William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.
- 700 Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face
701 model for pose and illumination invariant face recognition. In *2009 sixth IEEE international
conference on advanced video and signal based surveillance*, pp. 296–301. Ieee, 2009.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
models are unsupervised multitask learners. 2019.

- 702 Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: memory optimizations
703 toward training trillion parameter models. In *Proceedings of the International Conference for High*
704 *Performance Computing, Networking, Storage and Analysis*, SC '20. IEEE Press, 2020. ISBN
705 9781728199986.
- 706
707 Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System op-
708 timizations enable training deep learning models with over 100 billion parameters. In *Pro-*
709 *ceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery &*
710 *Data Mining*, KDD '20, pp. 3505–3506, New York, NY, USA, 2020. Association for Com-
711 puting Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3406703. URL <https://doi.org/10.1145/3394486.3406703>.
- 712
713 Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and
714 Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv preprint arXiv:2007.08501*,
715 2020.
- 716
717 Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael J Black. Learning to regress 3d face
718 shape and expression from an image without 3d supervision. In *Proceedings of the IEEE/CVF*
719 *Conference on Computer Vision and Pattern Recognition*, pp. 7763–7772, 2019.
- 720
721 Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to
722 instruction-follow them all. *arXiv preprint arXiv:2305.16355*, 2023.
- 723
724 Jianxin Sun, Qi Li, Weining Wang, Jian Zhao, and Zhenan Sun. Multi-caption text-to-face synthesis:
725 Dataset and algorithm. In *Proceedings of the 29th ACM International Conference on Multimedia*,
726 MM '21, pp. 2290–2298, New York, NY, USA, 2021. Association for Computing Machinery.
727 ISBN 9781450386517. doi: 10.1145/3474085.3475391. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3474085.3475391)
728 [3474085.3475391](https://doi.org/10.1145/3474085.3475391).
- 729
730 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
731 Liang, and Tatsunori B Hashimoto. Alpaca: A strong, replicable instruction-following model.
732 *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7, 2023.
- 733
734 Ayush Tewari, Michael Zollöfer, Hyeongwoo Kim, Pablo Garrido, Florian Bernard, Patrick Perez, and
735 Christian Theobalt. MoFA: Model-based Deep Convolutional Face Autoencoder for Unsupervised
736 Monocular Reconstruction. In *The IEEE International Conference on Computer Vision (ICCV)*,
737 2017.
- 738
739 Ayush Tewari, Michael Zollhoefer, Florian Bernard, Pablo Garrido, Hyeongwoo Kim, Patrick Perez,
740 and Christian Theobalt. High-fidelity monocular face reconstruction based on an unsupervised
741 model-based face autoencoder. *IEEE transactions on pattern analysis and machine intelligence*,
742 42(2):357–370, 2018a.
- 743
744 Ayush Tewari, Michael Zollhöfer, Pablo Garrido, Florian Bernard, Hyeongwoo Kim, Patrick Pérez,
745 and Christian Theobalt. Self-supervised multi-level face model learning for monocular recon-
746 struction at over 250 hz. In *Proceedings of the IEEE conference on computer vision and pattern*
747 *recognition*, pp. 2549–2559, 2018b.
- 748
749 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
750 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
751 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 752
753 M. Wu, H. Zhu, L. Huang, Y. Zhuang, Y. Lu, and X. Cao. High-fidelity 3d face generation from natural
754 language descriptions. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
755 *(CVPR)*, pp. 4521–4530, Los Alamitos, CA, USA, jun 2023a. IEEE Computer Society. doi:
10.1109/CVPR52729.2023.00439. URL [https://doi.ieeecomputersociety.org/](https://doi.ieeecomputersociety.org/10.1109/CVPR52729.2023.00439)
10.1109/CVPR52729.2023.00439.
- 756
757 Menghua Wu, Hao Zhu, Linjia Huang, Yiyu Zhuang, Yuanxun Lu, and Xun Cao. High-fidelity 3d
758 face generation from natural language descriptions. In *Proceedings of the IEEE/CVF Conference*
759 *on Computer Vision and Pattern Recognition*, pp. 4521–4530, 2023b.

756 Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal
757 llm. *arXiv preprint arXiv:2309.05519*, 2023c.
758

759 Cuican Yu, Guansong Lu, Yihan Zeng, Jian Sun, Xiaodan Liang, Huibin Li, Zongben Xu, Songcen Xu,
760 Wei Zhang, and Hang Xu. Towards high-fidelity text-guided 3d face generation and manipulation
761 using only images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
762 pp. 15326–15337, 2023.

763 Longwen Zhang, Qiwei Qiu, Hongyang Lin, Qixuan Zhang, Cheng Shi, Wei Yang, Ye Shi, Sibe
764 Yang, Lan Xu, and Jingyi Yu. Dreamface: Progressive generation of animatable 3d faces under
765 text guidance. *arXiv preprint arXiv:2304.03117*, 2023.
766

767 Ding Zheng, Zhang Cecilia, Xia Zhihao, Jebe Lars, Tu Zhuowen, and Zhang Xiuming. Diffusionrig:
768 Learning personalized priors for facial appearance editing. In *Proceedings of the IEEE/CVF
769 Conference on Computer Vision and Pattern Recognition*, 2023.

770 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: En-
771 hancing vision-language understanding with advanced large language models. *arXiv preprint
772 arXiv:2304.10592*, 2023.
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810 A APPENDIX

811
812 A.1 SPECULATIVE FACE GENERATION BENCHMARK

813
814
815 As an AI visual assistant specializing in human face analysis, your task is to infer possible
816 activities, events, and emotional states a person might be experiencing based on a detailed
817 visual and textual description of their face. Your goal is to focus on the high-level emotional
818 context and plausible scenarios that this person could be engaging in, rather than anatomical
819 details. For each facial description, consider the following key questions before formulating
820 your response:

- 821
822 1. What emotion is the person likely experiencing? How does this emotion differ from a
823 typical representation of this feeling?
824 2. What specific activity might this individual be participating in, based on their emotional
825 state and facial expression?
826 3. What event could have triggered the current emotional expression or facial state?
827 4. Could the individual be engaged in other parallel activities that are influencing their
828 expression?

829
830 Once you've considered these questions, craft 5 distinct facial descriptions, each beginning
831 with "This person," followed by one or two sentences that clearly suggest a plausible activity
832 or situation. Ensure that each description provides a rich context that allows the user to
833 imagine and even replicate the facial expression if they were in that situation. Avoid vague or
834 general terms, and be as diverse as possible in your interpretations.

835 Example answers and face descriptions:

836 Answer to the questions:

- 837 1. The individual appears to be deeply worried, but there's an undertone of surprise not typical
838 of worry alone.
839 2. The individual could be reviewing a critical work email and is taken aback by unexpected
840 information.
841 3. It seems like the individual just encountered an unexpected delay for an important event.
842 4. Additionally, the person might be contemplating a difficult decision while trying to process
843 sudden news.

844 Facial descriptions:

- 845 1. This person looks like they have just received news that their flight has been canceled at
846 the last minute.
847 2. This person is struggling to concentrate during a crucial work presentation, trying to mask
848 their frustration.
849 3. This person might be reading an intense plot twist in a book, causing both confusion and
850 intrigue.
851 4. This person seems to be on the verge of delivering uncomfortable feedback to a colleague.
852 5. This person appears to be caught off guard in a meeting, unexpectedly asked to answer a
853 difficult question.

854 Figure 5: Prompts for querying GPT-4V to convert explicit text descriptions and facial images into
855 implicit descriptions

856
857 Due to the lack of public datasets and benchmarks that provide implicit facial texts paired with
858 corresponding 3D faces, , we choose to construct the Speculative Face Generation Benchmark with
859 the help of powerful GPT4-V OpenAI et al. (2024). Inspired by Feng et al. (2024), we extract face
860 images and human annotated facial descriptions to GPT4-V to generate implicit facial descriptions.
861 Specifically, we select 74 high-quality facial images from CelebAHQ dataset Karras et al. (2018) and
862 their explicit descriptions from CelebAText dataset Sun et al. (2021). Using the prompt detailed in
863 Figure 5, we instruct the GPT4-V to analyze our provided images and detailed human descriptions
and generate appropriate implicit descriptions to reflect the possible activities and the emotion

864 states associated with the faces. For each case, we ask GPT4-V to produce five potential implicit
 865 descriptions. We then fed these five candidate descriptions, along with the corresponding facial image,
 866 back into GPT-4V and requested it to select the implicit description that best matches the image. This
 867 process yielded a total of 444 responses from GPT-4V and 74 candidate text-face pairs. After manual
 868 verification, we finalized 64 test cases for our Speculative Face Generation benchmark.

870 A.2 IMPLEMENTATION DETAILS ON LLaVA-KEY BASELINE

871 For Image-based 3D Face Reconstruction and text-based 3D Face Reconstruction, we develop
 872 a baseline based on VLM called LLaVA-Key where the human face is represented as text.
 873 Following Feng et al. (2024), we represent the human face through the 68 landmarks on 2D
 874 images as landmarks tracks the locations of eyes, nose, mouth and so on, which also contains rich
 875 information about faces. We utilize templates like "USER: <IMAGE> Please estimate
 876 the 68 facial landmarks coordinates. The output format should be
 877 Jawline-1:(x1,y1),Jawline-2:(x2,y2),... ASSISTANT: The detected
 878 landmarks are Jawline-0:(41, 88),Jawline-1:(43, 107), ...,Inner
 879 Lip-67:(101, 159)." for image-to-face data and templates like "USER: There
 880 is a person with the following description:{description} Please
 881 estimate the 68 facial landmarks coordinates. The output format
 882 should be Jawline-1:(x1,y1),Jawline-2:(x2,y2),... ASSISTANT: The
 883 detected landmarks are Jawline-0:(41, 88),Jawline-1:(43, 107),
 884 ...,Inner Lip-67:(101, 159)." for text-to-face data where {description} would
 885 be replaces with text description for faces generated by LLaVA. And we use an off-the-shelf face
 886 detector to provide the coordinates landmarks for each face. LLaVA model is finetuned to produce
 887 2D coordinates of 68 facial landmarks with the formatted data. After finetuning, LLaVA would be
 888 prompted to predict the landmarks in the test set and an optimization-based method is used to fit
 889 3DMM for the estimated landmarks. For image-to-face data, we utilize $loss_{face}$ defined in 3 to
 890 optimize the 3DMM parameters on estimated landmarks and input image. For text-to-face data,
 891 we only utilize the landmark loss and regularization defined in 3.2 for optimization. Experiments
 892 reflect that representing human face as a new modality outperform naively encoding human face in
 893 language.

894 A.3 FACE-CENTRIC CONVERSATION GENERATION

895 To enrich the diversity of question-answer pairs for image-to-face data, we inquiry LLaVA with the
 896 face image and the questions listed in 6 and collect the answers from LLaVA to construct face-centric
 897 conversations for each face image. During training, we would randomly pick question-answer pairs
 898 from these generated conversation and fuse the conversation with task-static template presented in
 899 4.1 as the text of training data.

901 Table 6: The list of instructions for constructing face-centric conversations.

-
- 904 • "How is the person’s hair styled?"
 - 905 • "What colors dominate this image?"
 - 906 • "Based on the attire and styling, can you infer anything about the event or occasion for
 907 this photo?"
 - 908 • "Can you describe the person’s expression?"
 - 909 • "Is there any indication of where this person might be?"
 - 910 • "What is the person wearing?"
-

915 A.4 TEXT-TO-FACE DATA GENERATION

916 As there is no public text-to-3DMM data available, we propose a way to utilize the powerful pretrained
 917 VLM to build a connection between face image with textual description. Specifically, we design

a template to inquiry pretrained LLaVA to output detailed description about human faces, which is presented as following:

Analyze the image and generate a detailed textual description of the human face it contains. Focus on the following aspects:

1. Face Shape: Description of the jawline shape (e.g., square, round, oval, heart-shaped). Forehead size and shape (e.g., wide, narrow, rounded). Cheekbone structure (e.g., high, low, prominent).
 2. Face Expression: Eyebrows (e.g., arched, straight, furrowed). Eyes (e.g., wide open, squinting, normal). Mouth (e.g., smiling, frowning, neutral). Additional details if any specific expression is featured (e.g., wrinkling of forehead, dimples).
 3. Face Color: Skin tone (e.g., fair, olive, dark, light). Any distinct color features such as freckles, rosiness, tan lines. Makeup if applicable (e.g., lipstick shade, eyeshadow color).
 4. Face Lighting: Direction of the light source (e.g., frontal, side, backlit). Intensity of the light (e.g., soft, harsh, moderate). Shadows observed on the face (specify areas such as under eyes, neck).
 5. Pose of Head: Mention the orientation of the head (e.g., facing forward, tilted to the side, looking upwards).
- Please give a response starting with 'He' or 'She'.

Figure 6: Prompts for querying LLaVA to generate explicit text descriptions on face images

The question covers many requirement on detailed description for many facial attributes and the response from pretrained LLaVA would be collected as the description for the face image and the 3DMM parameter of the corresponding face. Then we utilize these generated description with unsupervised face loss described in 3 to guide the model optimization.

A.5 EVALUATION METRICS

For text-based face generation tasks, we generally use three metrics to measure the quality of the predicted 3DMM: Chamfer Distance(CD), Complete Rate(CR), and Relative Face Recognition Rate(RFRR). We will elaborate about how to compute each metric below:

- Chamfer Distance(CD): CD measures the similarity between two sets of point clouds. Given the predicted 3D mesh M_p and the ground truth 3D mesh M_g , Chamfer Distance is defined as:

$$CD(M_p, M_g) = \sum_{a \in M_p} \min_{b \in M_g} \|a - b\|_2 + \sum_{b \in M_g} \min_{a \in M_p} \|a - b\|_2, \quad (6)$$

- Complete Rate(CR): CR is used to evaluate the completeness of a 3D point cloud, which can be formulated as the ratio between matched points and all points. The mathematical formulation of CR can be defined as:

$$CR = \frac{P_0}{P} \quad (7)$$

P_0 is the number of points with CD value less than 10mm and P is the number of all points.

- Relative Face Recognition Rate(RFRR): RFRR measures how well the reconstructed 3D face preserves identity information in comparison to the original face. For our case, we follow Wu et al. (2023a) and render the predicted mesh and the ground truth with the same pose and the lighting. Then we compute the cosine similarity in the feature space of ArcFace Deng et al. (2019a) to measure the identity preservation.

For image-based face reconstruction task, we measure the L2 photometric error and the L2 landmark error by comparing the rendered 3D face and the ground truth 2D images as they can directly reflect the quality of geometry and texture for the estimates meshes.

Table 7: Influence of LLM towards human face understanding

| Method | photo | keypoint |
|----------------|-------|----------|
| CLIP ViT + MLP | 0.133 | 7.6px |
| FaceGPT | 0.103 | 3.0px |

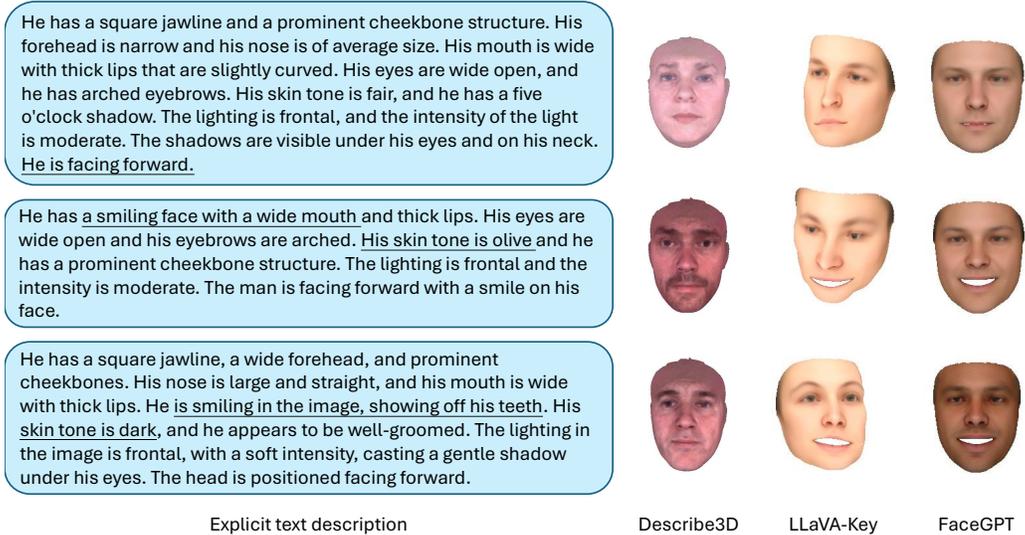


Figure 7: We present some examples for explicit text-based face generation tasks. While some methods may overlook specific underlined details in the descriptions, FaceGPT consistently produces 3D faces with superior semantic alignment to the provided text.

- Photometric error: Photometric error can measure the quality of geometry and texture for 3D face estimation jointly. The error E_{photo} can be formulated as:

$$E_{photo} = \|A \odot (x_{img} - \hat{y}_{rec})\|_2 \tag{8}$$

A is the skin region in the ground truth face image.

- Landmark error: Landmark error can measure the quality of the geometry of 3D face and the accuracy of the estimated camera pose. The error E_{lm} can be formulated as:

$$E_{lm} = \|LM_{img} - LM_{rec}\|_2 \tag{9}$$

A.6 MORE RESULTS

A.6.1 INFLUENCE OF LLM ON FACE UNDERSTANDING

Classical methods mainly rely on a vision encoder for regressing the 3DMM face parameters. In contrast, LLaVA utilizes a combination of frozen vision encoder and an LLM to perceive information from images. To study the effect of the LLM head, we train a baseline model with the vision encoder of LLaVA and an MLP to predict the human face parameters directly. The results are presented in Table 7 and show that the additional LLM helps in predicting better face parameters based on the visual representations from the frozen encoder.

A.6.2 COMPARISON BETWEEN SUPERVISED FACE LOSS AND UNSUPERVISED FACE LOSS

Recent VLM-based works on segmentation and pose estimation generally rely on supervised learning with ground truth data or pseudo ground truth data. To study the influence of supervision on the model’s capability, we compare the model trained with supervised 3DMM losses and self-supervised

Table 8: Comparison when learning with supervised and unsupervised losses.

| Method | photo. | keyp. |
|----------------|--------|-------|
| FaceGPT(sup) | 0.092 | 2.1px |
| FaceGPT(unsup) | 0.103 | 3.0px |

2D face losses. For the supervised 3DMM loss, we extract the 3DMM parameters on the same CelebA-HQ trainset with the state-of-the-art face reconstruction method FOCUS (Li et al., 2023). L1 loss is used for optimization on 3DMM ground truth. We observe in Table 8, that the model utilizing a supervised 3DMM loss outperforms its unsupervised counterpart, indicating that the performance ceiling is not reached yet and improvements on the self-supervised training could potentially lead to further performance gains.

A.6.3 QUANTITATIVE RESULTS FOR EXPLICIT TEXT-BASED FACE GENERATION

We present visual results for explicit text-based face generation in Figure 7. Compared to other methods, the output of FaceGPT are generally better aligned to the text descriptions.