



CAIRO UNIVERSITY

FACULTY OF ENGINEERING

DEPARTMENT OF COMPUTER ENGINEERING

Retratista



A Graduation Project Report Submitted
to
Faculty of Engineering, Cairo University
in Partial Fulfillment of the requirements of the degree
of
Bachelor of Science in Computer Engineering.

Presented by

Mohamed Shawky Zaky AbdelAal Sabae

Remonda Talaat Eskarous

Mohamed Ahmed Mohamed Ahmed

Mohamed Ramzy Helmy Ibrahim

Supervised by

Dr. Mayada Hadhoud

26th July, 2021

All rights reserved. This report may not be reproduced in whole or in part, by photocopying or other means, without the permission of the authors/department.

Abstract

المخلص

Acknowledgement

Contents

Abstract (English)	1
Abstract (Arabic)	2
Acknowledgement	3
Table of Contents	6
List of Figures	7
List of Tables	8
List of Abbreviation	9
List of Symbols	10
Contacts	11
1 Introduction	13
1.1 Motivation and Justification	13
1.2 Project Objectives and Problem Definition	13
1.3 Project Outcomes	13
1.4 Document Organization	13
2 Market Feasibility Study	14
2.1 Targeted Customers	14
2.2 Market Survey	14
2.2.1 FaceAPP	14
2.2.2 PicsArt	14
2.2.3 Facetune2	14
2.2.4 Booth Apps	14
2.3 Business Case and Financial Analysis	14
3 Literature Survey	15
3.1 Generative Models	15
3.2 Face Modelling and Generation	15
3.3 Comparative Study of Previous Work	15
3.4 Implemented Approach	15
4 System Design and Architecture	16
4.1 Overview and Assumptions	16
4.2 System Architecture	18
4.2.1 Block Diagram	19
4.3 Module 1 : Speech Recognition	20

4.3.1	Functional Description	20
4.3.2	Modular Decomposition	20
4.3.3	Design Constraints	20
4.3.4	Other Description	20
4.4	Module 2 : Text Processing	21
4.4.1	Functional Description	21
4.4.2	Modular Decomposition	21
4.4.3	Design Constraints	21
4.4.4	Other Description	21
4.5	Module 3 : Face Code Generation	22
4.5.1	Functional Description	22
4.5.2	Modular Decomposition	22
4.5.2.1	Feature Directions Generation	23
4.5.2.2	Initial Seed Generation	25
4.5.2.3	Latent Manipulation	25
4.5.3	Design Constraints	26
4.5.4	Synthetic Image Clustering	27
4.6	Module 4 : Code-to-Face Translation	28
4.6.1	Functional Description	28
4.6.2	Modular Decomposition	28
4.6.3	Design Constraints	30
4.7	Module 5 : Face Refinement	31
4.7.1	Functional Description	31
4.7.2	Modular Decomposition	31
4.7.3	Design Constraints	32
4.8	Module 6 : Multiple Head Poses Generation	33
4.8.1	Functional Description	33
4.8.2	Modular Decomposition	33
4.8.3	Design Constraints	33
4.8.4	Other Description	33
4.9	Module 7 : Web Application	34
4.9.1	Functional Description	34
4.9.2	Modular Decomposition	34
4.9.3	Design Constraints	34
4.9.4	Other Description	34
4.10	Other Approaches	35
5	System Testing and Verification	36
5.1	Testing Setup	36
5.2	Testing Plan and Strategy	36
5.2.1	Module Testing	37
5.2.1.1	Speech Recognition	37
5.2.1.2	Text Processing	37
5.2.1.3	Code Generation	37

5.2.1.4	Code-to-Face Translation	38
5.2.1.5	Face Refinement	39
5.2.1.6	Multiple Head Poses Generation	41
5.2.1.7	Web Application	41
5.2.2	Integration Testing	41
5.3	Testing Schedule	44
5.4	Comparative Results to Previous Work	44
6	Conclusions and Future Work	46
6.1	Faced Challenges	46
6.2	Gained Experience	46
6.3	Conclusions	46
6.4	Future Work	46
A	Development Platforms and Tools	48
A.1	Hardware Platforms	48
A.2	Software Tools	48
B	Use Cases	48
C	User Guide	48
D	Code Documentation	48
E	Feasibility Study	48

List of Figures

4.1	Block diagram of complete system architecture	19
4.2	Block diagram of application design	19
4.3	Detailed block diagram of the three core modules workflow	22
4.4	Illustration of feature directions in latent space	23
4.5	Illustration of orthogonalization relative to a reference vector	24
4.6	Illustration of directions scale using age direction	26
4.7	Style-based GAN architecture against traditional GAN	28
4.8	Illustration of sequential navigation and invertibility in a $2D$ latent space	30
5.1	The results of moving along some extracted feature directions.	38
5.2	An example of the consistency in reaching the required facial attributes starting from initial random vector.	39
5.3	The results of sequential navigation along certain feature directions.	40
5.4	Samples of correctly generated face portrait from textual description.	42
5.5	Samples of correctly generated face portrait from textual description.	43
5.6	Samples of incorrectly generated face portrait from textual description.	44
5.7	Plot of FID score of different pipelines against different number of test images (lower is better).	45

List of Tables

5.1	Angles between different feature directions using a subset of the considered facial features (closer to 90 degrees is better).	37
5.2	LPIPS values against the number of navigated directions for sample text (lower is better).	39
5.3	FID scores comparison on different number of test images (lower is better). . .	45
5.4	LPIPS comparison with Faces à la Carte (lower is better).	45
5.5	The execution time of different stages of the core of our system (measured in seconds).	46

List of Abbreviation

List of Symbols

Contacts

1 Introduction

1.1 Motivation and Justification

1.2 Project Objectives and Problem Definition

1.3 Project Outcomes

1.4 Document Organization

2 Market Feasibility Study

2.1 Targeted Customers

2.2 Market Survey

2.2.1 FaceAPP

2.2.2 PicsArt

2.2.3 Facetune2

2.2.4 Booth Apps

2.3 Business Case and Financial Analysis

3 Literature Survey

3.1 Generative Models

3.2 Face Modelling and Generation

3.3 Comparative Study of Previous Work

3.4 Implemented Approach

4 System Design and Architecture

In this chapter, we discuss our working pipeline and system architecture in details. Generally, our system takes a speech note, textual description or numerical attributes as an input. It processes the input description and outputs the initial human face portrait that corresponds to the given description. Afterwards, the user is allowed to manually control some facial attributes and morphological features and to rotate the face and render it in multiple poses. In the first section, we give an overview about the final system. Then, we discuss the final system architecture in the second section. In the subsequent sections, each module implementation is discussed in details. In the last section, we discuss the other conducted experiments, why we choose this final system and suggestions that can possibly improve the other experiments.

4.1 Overview and Assumptions

As mentioned above, our system basically enables the user to describe a human face in words or using numerical values and turns it into a full human face portrait that can be manipulated and rendered in multiple poses. The system relies heavily on generative models and text processing, both are iteratively designed to obtain the required results. The overall flow can be described as follows :

- The input speech notes are translated to text.
- The textual description (extracted from speech input or manually entered) is processed to extract the numerical values of the required facial features.
- The numerical values are used generate a face embedding vector that encodes the facial attributes in low dimensional space ($512D$).
- A generative model is specifically designed to translate from the low dimensional embedding into the full face portrait (1024×1024).
- The generated face portrait can be further refined by navigating the face embedding space and re-generating the face portrait.
- Once the user settles on the final face portrait, the system can render that face in multiple poses to provide further identification.

The previous flow provides a very versatile framework to generate face portrait and adjust it to your liking. However, there is an extremely large number of facial attributes and morphological features to describe a human face. Consequently, we have to choose a descriptive subset of these attributes to consider in the face description. We consider 32 facial attributes for face description, which are listed as follows :

- Overall face :
 - Gender : Male / Female.
 - Age : Young / Old.

- Thickness : Chubby / Slim.
 - Shape : Oval / Circular.
 - Skin Color : Black / White.
 - Cheeks : Normal / Rosy.
- Eyes :
 - Color : Black / Blue / Green / Brown.
 - Width : Wide / Narrow.
 - Eyebrows : Light / Bushy.
 - Bags Under Eyes : On / Off.
- Nose :
 - Size : Big / Small.
 - Pointy : On / Off.
- Ears :
 - Size : Big / Small.
- Jaw :
 - Mouth Size : Big / Small.
 - Lips Size : Big / Small.
 - Cheekbones : Low / High.
 - Double Chin : On / Off.
- Hair :
 - Color : Black / Blonde / Brown / Red / Gray.
 - Length : Tall / Short.
 - Style : Straight / Curly / Receding Hairline / Bald / with Bangs.
- Facial Hair :
 - Beard / None.
- Race :
 - White / Black / Asian.
- Accessories :
 - Glasses : Sight / Sun.
 - Makeup : On / Off.
 - Lipstick : On / Off.

4.2 System Architecture

Now, let's discuss our system architecture. The system consists of 6 modules, 3 core modules of the project and 3 auxiliary modules. These modules are deployed in a *web application* to provide an easy-to-use interface for face generation and manipulation. Figure 4.1 shows the complete block diagram of the system architecture. Meanwhile, figure 4.2 shows the application design and how the modules are deployed in a web application. The *core* modules are listed as follows :

- **Text Processing** : processes the input textual description and extracts the corresponding numerical values of facial attributes. This problem is similar to *multi-label text classification*, however the outputs are normalized scores of facial attributes, which are designed carefully to match the *face code generation* process.
- **Face Generation** :
 - **Code Generation** : converts the numerical attributes values to be low dimensional face embedding. This is the most *important* and *innovative* module of our system, because it glues the desired attributes scored with the latent space of the generative model (used to generate the face), resulting in more accurate quality outputs.
 - **Code-to-Face Translation** : translates the low dimensional face embedding into the actual face portrait. For this purpose, we use StyleGAN2, which is a *state-of-art latent-based generative model*, whose latent space can be manipulated easily to fit our needs.

Meanwhile, the *auxiliary* modules are listed as follows :

- **Speech Recognition** : translates the input speech to textual description.
- **Face Refinement** : uses the same generative model to manually refine the generated face portrait through navigating the latent space.
- **Multiple Head Poses Generation** : rotates the generated face portrait and renders it into multiple poses.

We discuss each module in more details in the subsequent sections. Also, these modules are organized into a web application for easier usage, as shown in Figure 4.2. The application is divided into :

- **Web (Frontend)** : which contains the user interface and, also, the *speech recognizer*. The speech recognizer is moved to the frontend to reduce the network communication overhead between the web application and the server, as transmitting text is easier than transmitting speech. Moreover, the speech recognizer doesn't require high computational power, so it can be embedded in the web application.
- **Server (Backend)** : which is separated into two servers. First server contains the *text processor* and the *generative model* and serves the requests of face generation and refinement. Second server contains the *pose generator* and serves the requests of face rotation.

The two servers can communicate with each other to exchange the generated face portraits through TCP sockets. Meanwhile, the web application communicates and sends requests to the servers through HTTP REST API.

4.2.1 Block Diagram

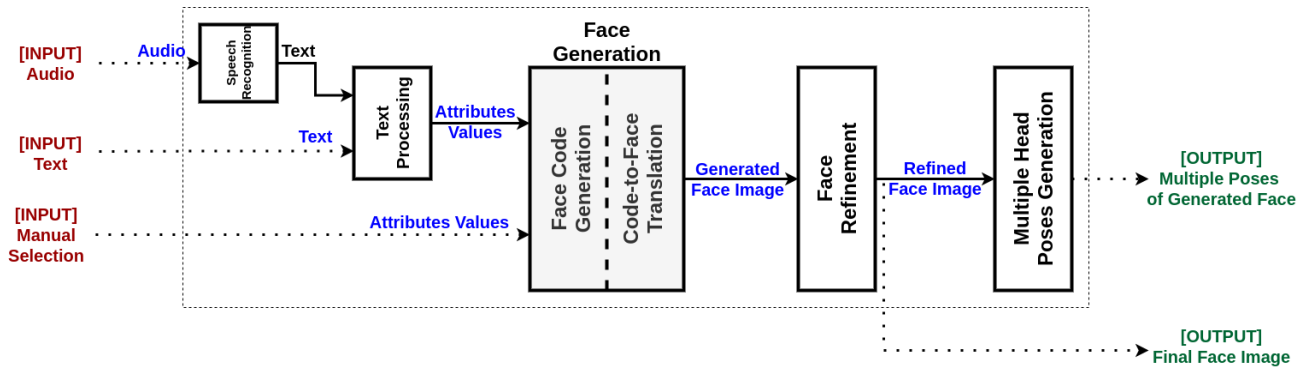


Figure 4.1: Block diagram of complete system architecture

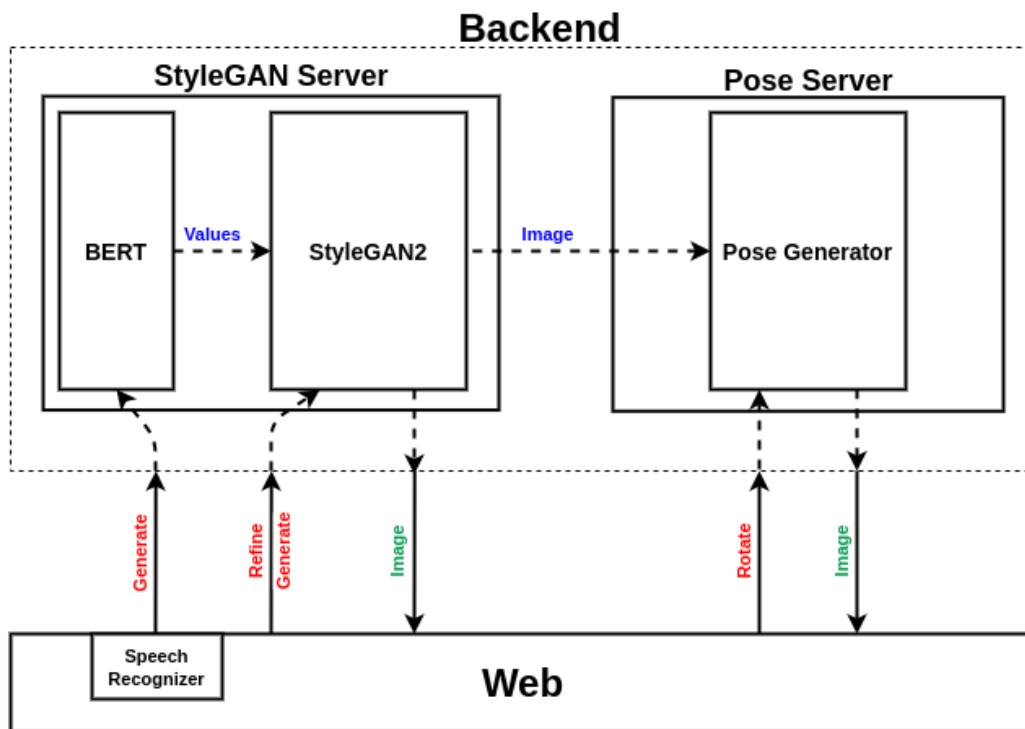


Figure 4.2: Block diagram of application design

4.3 Module 1 : Speech Recognition

4.3.1 Functional Description

4.3.2 Modular Decomposition

4.3.3 Design Constraints

4.3.4 Other Description

4.4 Module 2 : Text Processing

4.4.1 Functional Description

4.4.2 Modular Decomposition

4.4.3 Design Constraints

4.4.4 Other Description

4.5 Module 3 : Face Code Generation

Here, we discuss the face code generation from numerical values of facial attributes. This is the most important and innovative module in our system and the first stage of *face generation*. It's worth **noting** that we use both of the terms "*feature*" and "*attribute*" to refer to a facial attribute, like hair color or nose size.

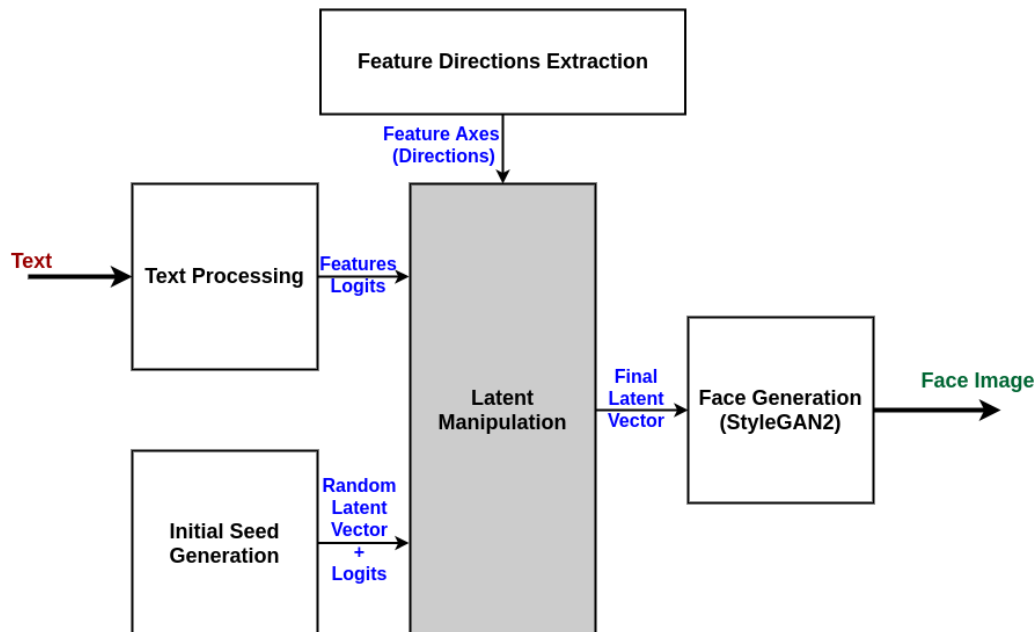


Figure 4.3: Detailed block diagram of the three core modules workflow

4.5.1 Functional Description

Figure 4.3 shows a block diagram of the interaction between the 3 core modules. We can see that the code generation module is the main driver of our face generation process. Generally, it converts the numerical attributes values (a.k.a. *logits*) into a face embedding vector that matches the design of the latent space of the face generator (*StyleGAN2*). Basically, it starts from an initial vector and uses the *required feature values* and *extracted feature directions* to transform this vector into the final latent vector, which is passed to the generative model.

- **Input :**
 - Numerical values of facial features (logits).
- **Output :**
 - Low dimensional face embedding vector (latent vector).

4.5.2 Modular Decomposition

As figure 4.3 tells, the code generation module can be torn down into 3 sub-modules, which are **latent manipulation**, **initial seed generation** and **feature directions extraction**. Each sub-module is discussed in details to show how they integrate to each other to achieve the desired goal.

4.5.2.1 Feature Directions Generation

Since, we use StyleGAN2 [1] as our generative model, we have a full $512D$ latent space that is used to encode the whole face attributes. The changes in this latent space maps to the generated face image and similar features occupies the same area in the latent space. Consequently, we have to come up with a way to extract the axes (*hyperplanes*) in this latent space to define each of our 32 facial features. These feature directions are, then, used to manipulate the latent vector, in order to map to the required face image.

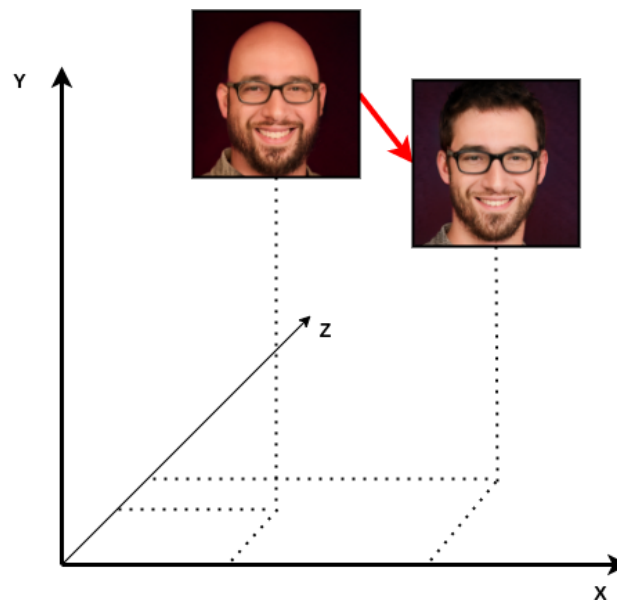


Figure 4.4: Illustration of feature directions in latent space

Figure 4.4 further illustrates the idea of feature directions in the latent space. Here, we plot two face images in a $3D$ latent space. We can see that the difference between the two images in the existence and the absence of the hair, thus the red arrow represents the *baldness* feature direction in that $3D$ latent space (moving along this particular vector causes hair density to change).

Our method of extracting the feature directions (*hyperplanes*) consists of 3 steps :

1. **Code-Image Pairs Generation and Classification** : First, we use StyleGAN2 to generate a large number of synthetic faces from random latent vectors. After so, we cluster the synthetic images (along with their latent vectors) according to each feature. The clustering can be based on discrete categories (like *hair color* or *race*) or continuous values (like *hair length* or *nose size*). We randomize the synthetic images in each clustering process to have better generalization and to cope with potential generation noise. For classification and regression, we use one of three possible methods, which are **manual labelling**, **classical image processing techniques** and **neural networks**. Thus, the output of this process is different groups of synthetic images sharing common facial features, along with their latent vectors.
2. **Feature Directions Fitting** : Now, we have a set of latent vectors (x) and their corresponding feature values (y). It's required to find a set of feature directions that satisfies

the mapping between feature vectors and values. This problem can be formulated as :

$$Y = A_f \cdot X \quad (1)$$

Where A_f is the axis (direction) of feature f .

We can obtain the solution to this equation in a closed form. However, due to the noise in both generation and classification, along with the non-linear nature of the problem, we opt to use *ML* methods, specifically **Logistic Regression** and **SVM** to get an *approximate solution*. Meanwhile, we cannot see any difference between the two methods, as they yield almost the same results.

Finally, the generated feature directions are normalized to unit vectors :

$$A_{unit} = \frac{A}{||A||} \quad (2)$$

3. **Directions Orthogonalization** : Facial features entanglement is one of the most difficult challenged of face generation. Some attributes in the human face tend to be extremely entangled by nature. For example, Asians rarely have curly hair, a woman cannot have beard and a man cannot put on makeup. Since StyleGAN2 is trained and tuned on **FFHQ** dataset [2], which contains real human faces, it is normal to notice some entanglement between some features. Consequently, the feature directions have to be further disentangled by using *orthogonalization*. The orthogonalization process is done iteratively, starting from the most accurate feature directions. We orthogonalize other feature directions on the accurate ones, so that we have completely independent feature directions, where tuning one direction doesn't affect the others. The directions are orthogonalized as follows :

$$A_{proj} = (A \cdot B_{unit})B_{unit} \quad (3)$$

$$A_{orthogonal} = A - A_{proj} \quad (4)$$

Figure 4.5 visually illustrates the *directions orthogonalization* process on 2D vectors.

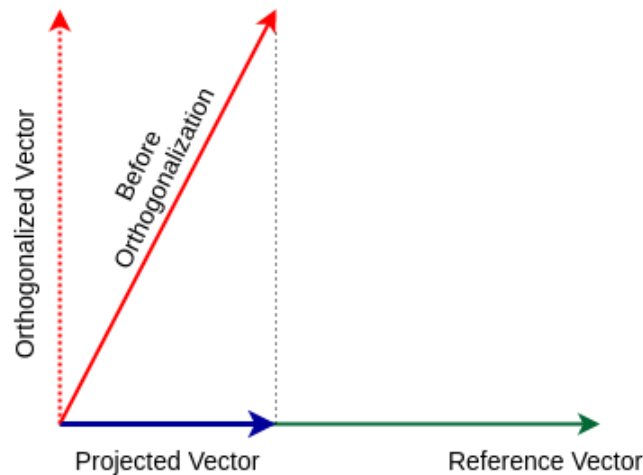


Figure 4.5: Illustration of orthogonalization relative to a reference vector

To ensure convergence to reasonable set of feature directions, we use a threshold margin to stop the orthogonalization process, which is from 85 to 95 degrees (5 degrees on each side of normal angle).

4.5.2.2 Initial Seed Generation

In order to avoid noise and discontinuities in the latent space, we generate an initial random latent vector. This is done by generating a random $512D$ vector and then passing it through the *mapping network* of `StyleGAN2`, which is not invertible. This initial vector is, then, manipulated by sequential navigation along each feature direction (axis) with certain amounts. To get these amounts, we should know the component of the initial vector along each feature direction. We do that by simply performing a dot product between the initial vector and the unit vector of each feature direction. Thus, we have an initial latent vector and the numerical attributes values, it presents.

4.5.2.3 Latent Manipulation

This sub-module ingests all the inputs and produces the required latent vector (*face embedding*) that describes all of the required facial attributes. The inputs to this latent manipulation sub-module are *initial random vector* along with its logits, *text logits* and *feature directions*. The latent manipulation, simply, wants to realize the following transformation on the *initial random vector* :

$$E_{final} = E_{initial} + (l_{text} - l_{rand})D \quad (5)$$

Where E_{final} is the final latent (*embedding*) vector of dimensions 1×512 , $E_{initial}$ is the initial random vector of dimensions 1×512 , l_{text} is the text logits vector of dimensions 1×32 (remember that we consider 32 facial features), l_{rand} is the logits vector of the initial random vector of dimensions 1×32 and D is the feature directions matrix of dimensions 32×512 . The transformation includes calculating the difference between the required logits and the random logits and, then, use this difference to move the initial random vector along the feature directions to reach the final latent vector.

It might seem straight forward to perform this transformation. Unfortunately, it's not feasible to perform the transformation using direct matrix multiplication, mainly due to heavy *entanglement* between direction vectors even after *orthogonalization*. Also, the latent space of `StyleGAN2` can be very noisy in certain regions, so transformations have to be done carefully.

Consequently, the processing in this sub-module is done iteratively as follows :

- Both random and text logits are scaled from 0 to 1, which cannot have significant effect, when navigating using unit directions. Consequently, the inputs logits are scaled with the directions scale, which is obtained empirically to be from -4 to 4 , as shown in figure 4.6.
- The next step is to get the *difference* between *text logits* and *random logits*, which is of dimensions 1×32 . We call that **differentiated logits**. It's worth noting here that the input text usually contains a *subset* of the facial features. Consequently, *not all* the text logits

are set to specific values. So, when doing the *differentiation*, we set the differentiated logits of the *unmentioned facial features* to 0. So, it can be summarized as follows :

$$l_{diff} = \begin{cases} l_{text} - l_{rand} & l_{text} \neq None \\ 0 & otherwise \end{cases} \quad (6)$$

- Loop over each direction in feature directions :
 - Multiply the *differentiated logit* corresponding to the *current feature* with its *direction*.
 - Add the product to the current latent vector (starting with the *initial random vector*). The following equation summarizes these steps :

$$E_{next} = E_{prev} + l_{diff}[j] * D[j] \quad (7)$$

- Finally, to ensure that every transformation is independent of the subsequent transformations and that they are applied sequentially, we re-compute the *produced latent vector logits* and re-differentiate it with the *text logits*. This is done on every iteration of the latent manipulation process.

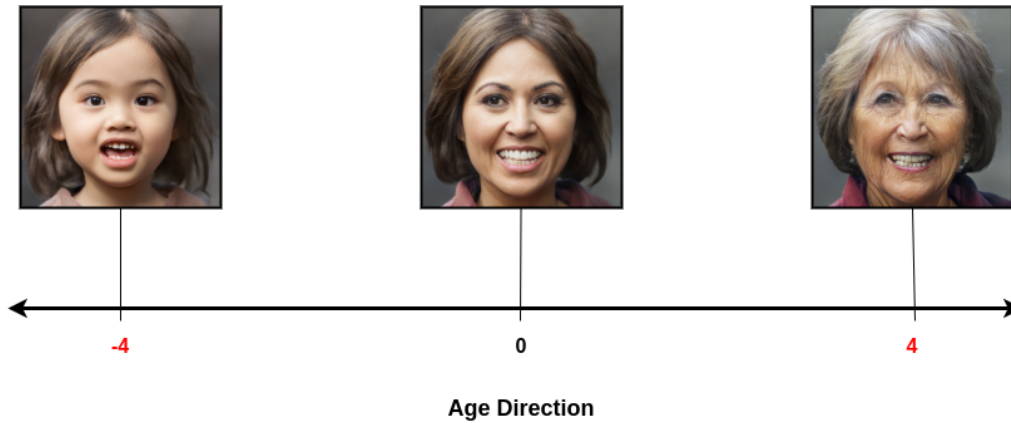


Figure 4.6: Illustration of directions scale using age direction

By applying the previous process, the *numerical value* of facial features extracted from text or manually entered by the user can be converted into a complete face embedding (*latent vector*) matching the required facial attributes. This vector can be passed to StyleGAN2 to translate it to a complete human face image.

4.5.3 Design Constraints

The design constraints of this module are enumerated as follows :

1. **Facial attributes entanglement** is the main challenge of the face code generation module. Naturally, human face attributes are related to each other. For example, Asians barely have curly hair, no woman cannot have beard and most women have long hair and wear makeup. We mentioned before that StyleGAN2 is trained and refined on FFHQ

dataset, which contains real human faces. Consequently, it's normal to see heavy entanglement between features in the latent space. Due to this *entanglement*, we have to perform extra computations to get decent results.

2. **Random initialization** of the latent vector can cause some issues with the final output, as the vector can be initialized in a noisy area of the latent vector. We try to *limit* the initialization of latent vector to a certain set of random vectors to avoid this effect. This method significantly reduces the *random initialization effect*, however it's not fully cured.
3. **Directions accuracy** can be a challenge as well. Some factors can negatively affect the feature directions accuracy. These factors include **synthetic image clustering** (whether *classification* or *regression*) and **directions fitting** process. We already discussed our solutions to this problem.

4.5.4 Synthetic Image Clustering

As we mentioned before, it's required to cluster the generated synthetic images, in order to be able to fit the feature directions. The *clustering* process can be performed through *classification* or *regression*. For example, features like hair and skin colors should be classified (grouped) into discrete categories, however features like hair length and mouth size should be assigned a continuous value. We do this task using one of three different methods :

1. **Manual labelling** is the first idea to come to our minds. It's straight forward to manually classify a group of faces according to a certain feature. This method is used with some features, however it's very cumbersome and only works for classification.
2. **Classical image processing techniques** are, also, used to classify images based on some features. Mainly, we use these techniques to detect *colors* like eye color and hair color. We use *morphological operators* and *classical segmentation* to detect *the eye* or *the hair* and then retrieve its color.
3. **Deep learning techniques** (*neural networks*) are used to perform regression on the rest of the features. We basically use *facial landmark detection* pretrained networks to detect the important facial landmarks, which is, then, used to calculate *sizes* and *distances*.

4.6 Module 4 : Code-to-Face Translation

Here, we discuss the process of converting the face embedding (*latent vector*) to a complete face image using *StyleGAN2* (our chosen generative model). We discuss our reasons for choosing this particular architecture and how we use it.

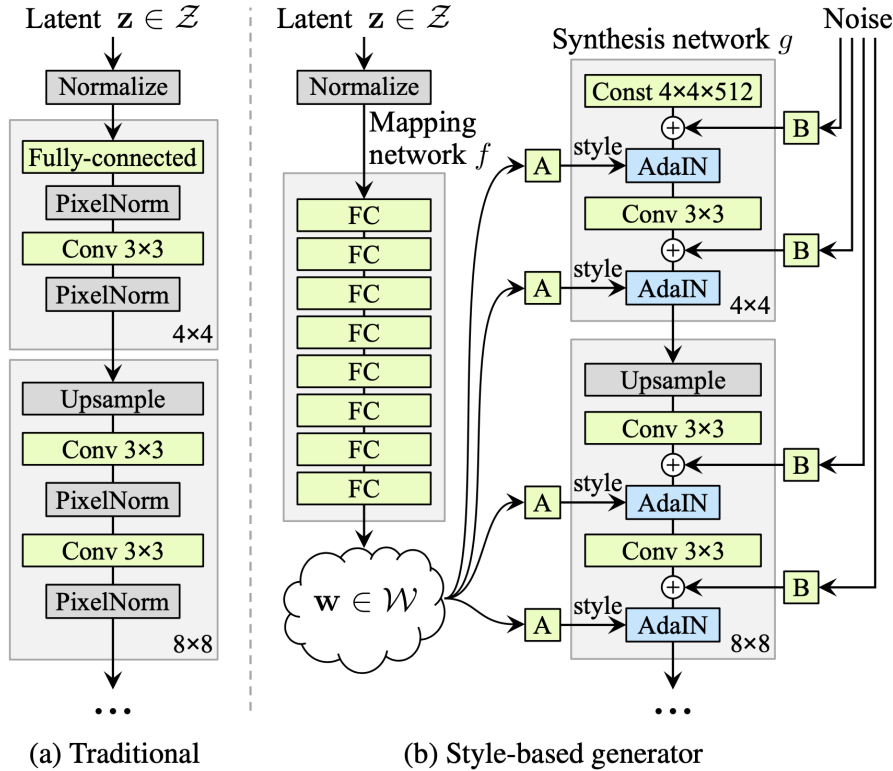


Figure 4.7: Style-based GAN architecture against traditional GAN

4.6.1 Functional Description

This module utilizes the power of *style-based* generative models, specifically *StyleGAN2* [1], to translate the required *latent vector* to a complete *human face image*. *Style-based GANs* (sometimes called *latent-based GANs*) can exert artistic control over the generated content (images, videos, text ..., etc). Consequently, we can tune it to fit our need and, iteratively, design it along with *code generation* 4.5 and *text processing* 4.4, in order to have a complete end-to-end pipeline for *text-to-face generation*.

- **Input :**
 - Low dimensional face embedding vector (latent vector).
- **Output :**
 - Complete human face image (portrait).

4.6.2 Modular Decomposition

As mentioned before, we opt to use *Style-based GANs* to be able to artistically control the output and designed the whole pipeline for *text-to-face generation*. Moreover, we choose

StyleGAN2, because it's one of the most popular and robust Style-based GANs in research literature. Also, it's relatively lightweight compared to other GANs used for the same purposes, but most importantly, StyleGAN2 excels at human face generation based on latent space. Figure 4.7 shows the original architecture of StyleGAN and how it is compared to traditional GANs. StyleGAN generator has two networks as follows :

- **Mapping network** creates nonlinear transformation to the input latent vector z ($512D$). This transformation is not invertible and results in a $512D$ latent vector w . This latent vector w is expanded into several $512D$ vector using affine transformation, which gives the *extended latent vector* $w+$. The extended latent vector $w+$ dimensions depend on the dimensions of the output image.
- **Synthesis network** generates the synthetic image from *normally-distributed noise* guided by the extended latent vector $w+$.

StyleGAN2 [1] is a newer version that follows the same architecture, but with some modifications to further improve the control over latent space and the quality of the outputs.

So, let's discuss how we adapted StyleGAN2 to our work :

- To provide a high fidelity results, we target 1024×1024 synthetic images. To achieve this, we have to use an extended latent vector $w+$ of dimensions 18×512 , meaning we repeat the latent vector w 18 times with *affine transformation* for each.
- To further improve feature directions disentanglement, we fine-tune StyleGAN2 using a subset of FFHQ dataset with increasing the weight of *perceptual path regularization* in the loss function. *Perceptual path regularization* in StyleGAN2 loss encourages the smooth mapping between latent and image spaces. So, when increasing it in certain directions, it highly penalizes the deviation between latent and image spaces in these directions giving more organized latent space. To avoid using the whole dataset, we use StyleGAN2 *adaptive discriminator augmentation (ADA)* [3] training methodology.
- Finally, we opt to remove the *mapping network* of StyleGAN2 generator and only use the *synthesis network*. This is mainly because :
 - The mapping network doesn't satisfy the *path length regularization*, so there is no smooth mapping between latent space z and image space (only with latent space w). This is discussed in the original paper [1] and our experiments support that.
 - The mapping network is not invertible, unlike the synthesis network. So, we cannot reverse the transformation from image space to latent space z . That's why we only work with latent space w to test the consistency of the results.
 - Removing the mapping network reduces the computations and the memory footprint, which is crucial in our case.

After the previous modifications, the output model can directly translate the face embedding vector (*latent space*) into a complete human face portrait.

Notice that using this methodology of **code-to-face translation** along with **code generation** makes the sequential navigation in the latent space is *easily invertible*, which eases the

generation and the refinement (discussed in 4.7) of the synthetic face and gives the system versatility and fault tolerance. Figure 4.8 illustrates the idea of *invertibility* of sequential edits on a $2D$ latent space example. Here, we show a simplified $2D$ latent space with 3 feature direction. AB vector represents the *hair color* direction and BC vector represents the *gender* direction. Point A starts with a *blonde girl*, moving along AB vector results in a *girl with black hair* at point B . Then, moving along BC vector gives us a *man with black hair* at point C . Consequently, moving along CA vector, which is the opposite of the resultant of AB and BC , gives us the original face, thus inverting the sequential changes.

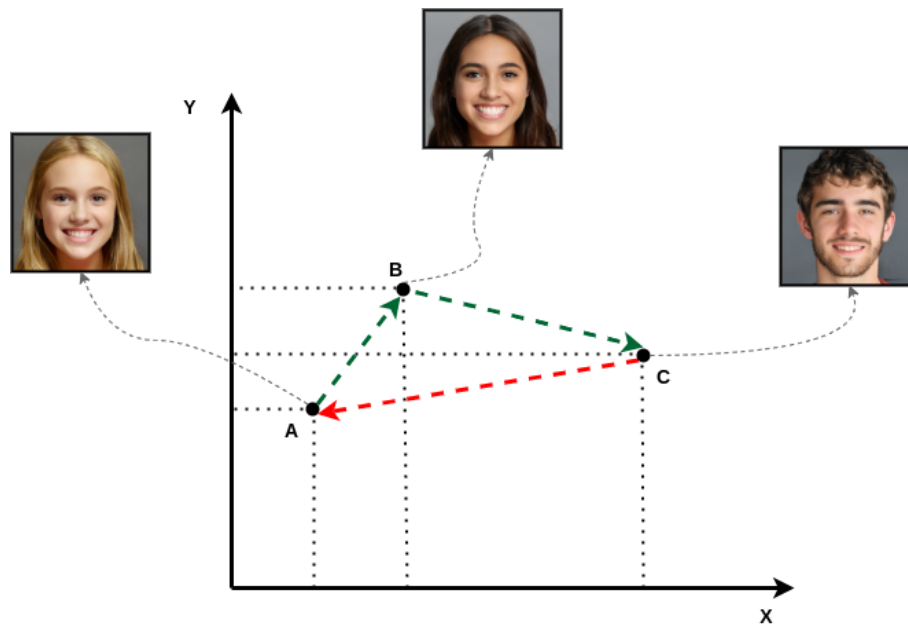


Figure 4.8: Illustration of sequential navigation and invertibility in a $2D$ latent space

4.6.3 Design Constraints

The basic design constraints for this module can be enumerated as follows :

1. **Network size** is surely one of our challenges. StyleGAN2 has a large memory footprint, as the case with most *deep generative models*. This constrains its training and deployment. We managed to remove the *mapping network*, which gives us a change to use the full *synthesis network*.
2. **Faces dataset** is, also, a constraint, because real images of human faces contains entanglement between facial features (*as discussed before*). So, we have to exert extra effort in the **code generation** module to solve some of this entanglement, emerging from real human faces datasets.

4.7 Module 5 : Face Refinement

Figure 4.8 shows that `StyleGAN2` latent space can be used to perform *directed manipulation* by using the feature directions extract in **code generation** 4.5. We discussed how we use this technique for *conditional sampling* from latent space by starting from an initial seed and editing it until we reach the desired latent vector. Now let's discuss our methodology of applying the same technique is *directed manipulation* of one or more features to be used in *face refinement*.

4.7.1 Functional Description

This module basically gives the users a chance to further refine the generated face portrait to their liking. The user can refine the generated face using the same facial features used in **face generation**, along with some additional features related to *face morphology*, which are hard to describe in words. We adapt the techniques used in *face generation* to perform *face refinement*.

- **Input :**
 - Generated human face image (portrait).
- **Output :**
 - Refined human face image (portrait).

4.7.2 Modular Decomposition

We utilize the same technique used *face generation* to perform *directed manipulation* of one or more features. This is exactly what is required by *face refinement*. Consequently, we opt to stick with our *face generation* pipeline, due to the following reasons :

- Latent manipulation in `StyleGAN2` yields better results than most attribute-editing GANs.
- We use the same network, which significantly reduces the memory footprint.
- Since `StyleGAN2` has an artistic control over its output, we can expand the number of target facial features, which improves the system scalability.

Using latent space navigation, we can easily perform *directed manipulation* on a subset of the facial features, which is required for *face refinement*. The *directed manipulation* over a single facial feature can be formulated as :

$$E_{refined} = E_{old} + \delta_{feature} * d_{feature} \quad (8)$$

Where $E_{refined}$ is the refined latent vector, E_{old} is the old latent vector, $\delta_{feature}$ is the change offset of the facial feature and $d_{feature}$ is the feature direction.

Subsequently, we pass the *refined latent vector* through the *synthesis network* of `StyleGAN2` to produce the new refined face portrait.

Finally, more feature directions are generated for the purpose of *face refinement*, which are hard to describe in words. These features are more related to *face morphology* and listed as follows :

- Distance between eyes.
- Distance between eyes and eyebrows.
- Distance between nose and mouth.
- Eyes opening.
- Mouth opening.
- Smiling or not.

These feature directions are obtained using the same technique described in 4.5, however all of these features are clustered using regression techniques. Overall, we use both facial features used in *face generation* and the new facial features (which sum up to 38 features) to refine the generated face.

4.7.3 Design Constraints

The only constraint, imposed on this module, is the difficulty to generate the feature directions. The new features describe spare facial attributes related to morphology, so it's hard to regress them. We use *facial landmark detection* techniques to detect the distances, in order to fit the feature directions. We couldn't further expand the morphological feature directions, due to *limited resources* and *time constraints*. Our *recommendation* is that using the facial landmarks detection can actually further expand the morphological features, so it's worth spending more time on it.

4.8 Module 6 : Multiple Head Poses Generation

4.8.1 Functional Description

4.8.2 Modular Decomposition

4.8.3 Design Constraints

4.8.4 Other Description

4.9 Module 7 : Web Application

4.9.1 Functional Description

4.9.2 Modular Decomposition

4.9.3 Design Constraints

4.9.4 Other Description

4.10 Other Approaches

5 System Testing and Verification

In this chapter, we discuss how the system is tested both on module level and integration level. We discuss our *performance metrics* and show the results of the system both *quantitatively* and *qualitatively*. We, also, show how our system performs compared to the baselines. Finally, we include the complexity analysis of our system for both time and memory.

5.1 Testing Setup

The testing environment of the whole system is basically targeting 4 main properties :

1. The quality of the output face images compared to the real human face images.
2. The ability of the system to capture the included facial features in the input description and how they actually map to the output.
3. The smooth and consistent mapping between the edits, imposed on the input description, and the changes of the output face image.
4. The Independence (*disentanglement*) between the different facial features.

We create our testing strategy to assess these 4 properties on the output face images. Also, we compare our results against the following baselines :

- StyleGAN2 [1] : We compare our results with the original *StyleGAN2* to check the output quality.
- Faces à la Carte [4] : This is the only previous research work that attempted *Text-to-Face Generation*.
- Image2StyleGAN [5] : Our feature directions extraction methodology is inspired by this work, which is based on the first version of *StyleGAN*.

5.2 Testing Plan and Strategy

To assess the 4 previously-mentioned properties, multiple metrics are used, which are listed as follows :

1. **Fréchet Inception Distance (FID)** [6] : This metric is an improvement over the traditional *inception score* to be able to measure the similarities between a set of real and synthetic images. Basically, *inception score* measures the ability of Inception V3 network [7] to classify a synthetic image into 1000 classes. However, *FID* measures the distance between synthetic and real images. This is done by extracting 2048D feature vector from each image using the Inception V3 network and then calculating the *Fréchet* distance using :

$$d^2 = ||\mathbf{u}_1 - \mathbf{u}_2||^2 + \text{Trace}(\mathbf{C}_1 + \mathbf{C}_2 - 2 * \sqrt{\mathbf{C}_1 * \mathbf{C}_2}) \quad (9)$$

Where \mathbf{u} is the feature-wise mean vector and \mathbf{C} is the covariance matrix of the feature vector.

2. **Learned Perceptual Image Patch Similarity (LPIPS)** [8] : This metric measures the smoothness of the mapping between the latent space edits and the output image changes. This metric takes as an input, two synthetic images. It uses a pretrained neural network to project them to a latent space. Then, it calculates the difference between the two latent vectors, along with the perceptual distance between the two images. Finally, it uses the two distances to calculate the final score. This metric is used in *StyleGAN2* paper to assess the *perceptual path length*.
3. **Edit Consistency Score** : We use this metric to ensure that the final facial attributes values are consistent with the input values after the *latent manipulation* process. This metric is simply calculated by projecting the final latent vector over all feature directions and compare it to the input values.
4. **Directions Disentanglement Score** : The disentanglement between feature directions are assessed by using the *angles* between each pairs of directions. Angles of values 85 to 95 degrees usually indicates low entanglement.

5.2.1 Module Testing

5.2.1.1 Speech Recognition

5.2.1.2 Text Processing

5.2.1.3 Code Generation

To test the quality of the generated feature directions that are used for *code generation*. We use two methods, which are **directions disentanglement scores** and **visual result** of moving along directions.

Table 5.1 shows the angles between the directions of a subset of features. Remember that angles in range 85 to 95 degrees indicate low entanglement. Consequently, we can infer that the number, provided by the table, are reasonable. For example, the angle between *gray hair* and *age* directions is 79.6 degrees, because old people normally have gray hair. Also, men are *not* likely to wear makeup, so the angle between *makeup* and *gender* directions is 107.7 degrees, same for *beard* with men.

Moreover, figure 5.1 shows the *visual results* of moving along some feature directions. We use these results to qualitatively measure the accuracy of the extracted feature directions.

Angles	Age	Gender	Beard	Gray Hair
Age	0.0	92.4	85.8	79.6
Gender	92.4	0.0	80.0	88.6
Makeup	88.0	107.7	100.5	94.5
Hair Length	89.7	95.9	90.6	96.6

Table 5.1: Angles between different feature directions using a subset of the considered facial features (closer to 90 degrees is better).

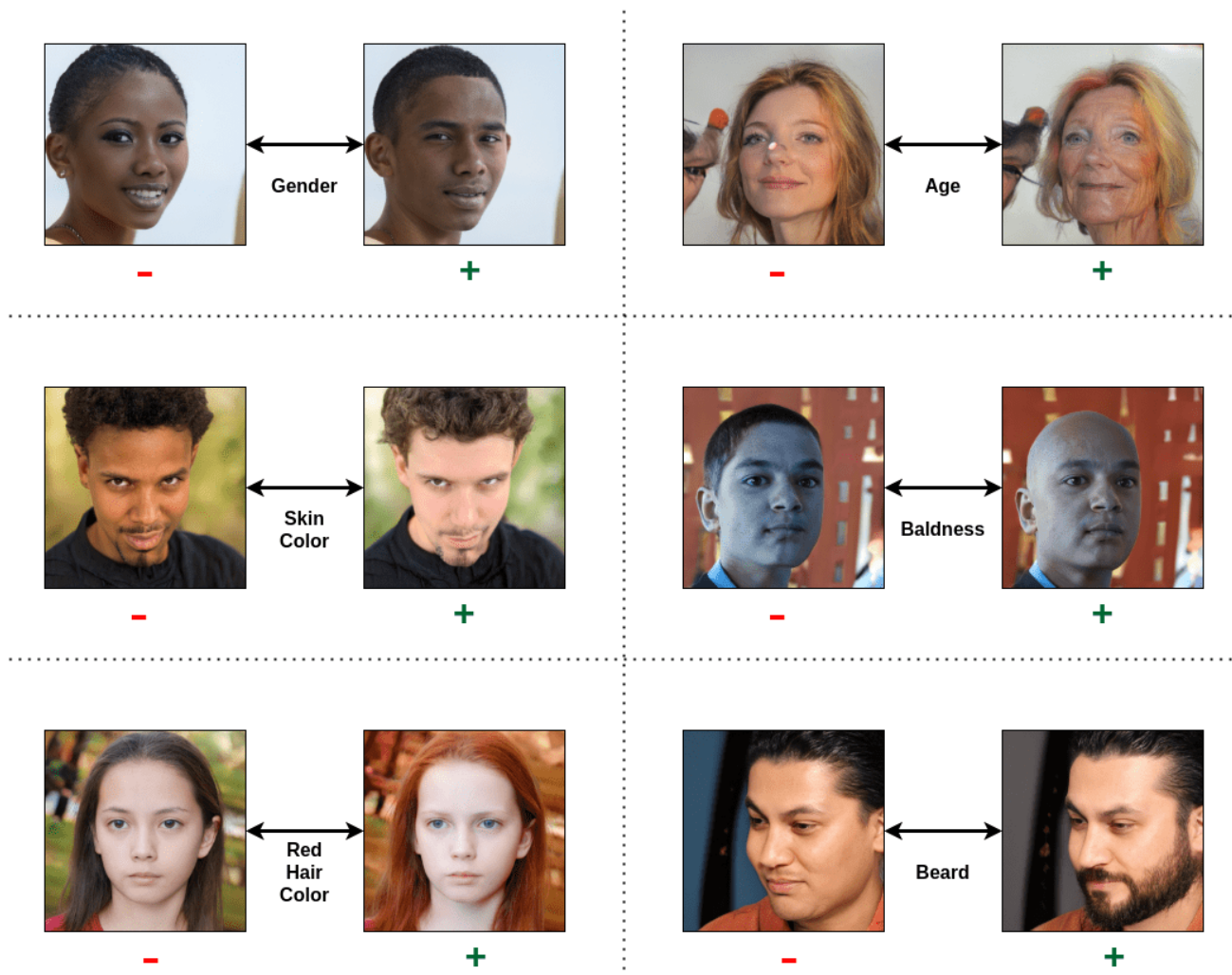


Figure 5.1: The results of moving along some extracted feature directions.

5.2.1.4 Code-to-Face Translation

This module is basically tested using integration with **code generation**, which is shown in the **integration** test. However, we perform testing on this module separately using **edit consistency score** of the facial attributes values of the output image and the input values.

As figure 5.2 shows, the system can convert a random vector (*on the left*) to the final latent vector (*on the right*) driven by the input values (*on top*). Also, we can see the consistency between the required values and the values corresponding to the generated face.

Also, table 5.2 shows the relation between *LPIPS* score (between the initial and output images) and the number of directions, navigated during face generation. We can see that the score increases, as the number of navigated directions increases. This is mainly because the output face image is far from the initial face image. However, we can see that some *anomalies* can occur, like the case of the input text *"Woman with lipstick and rosy cheeks"*, which navigates along only 3 directions, but gives a high *LPIPS* score.

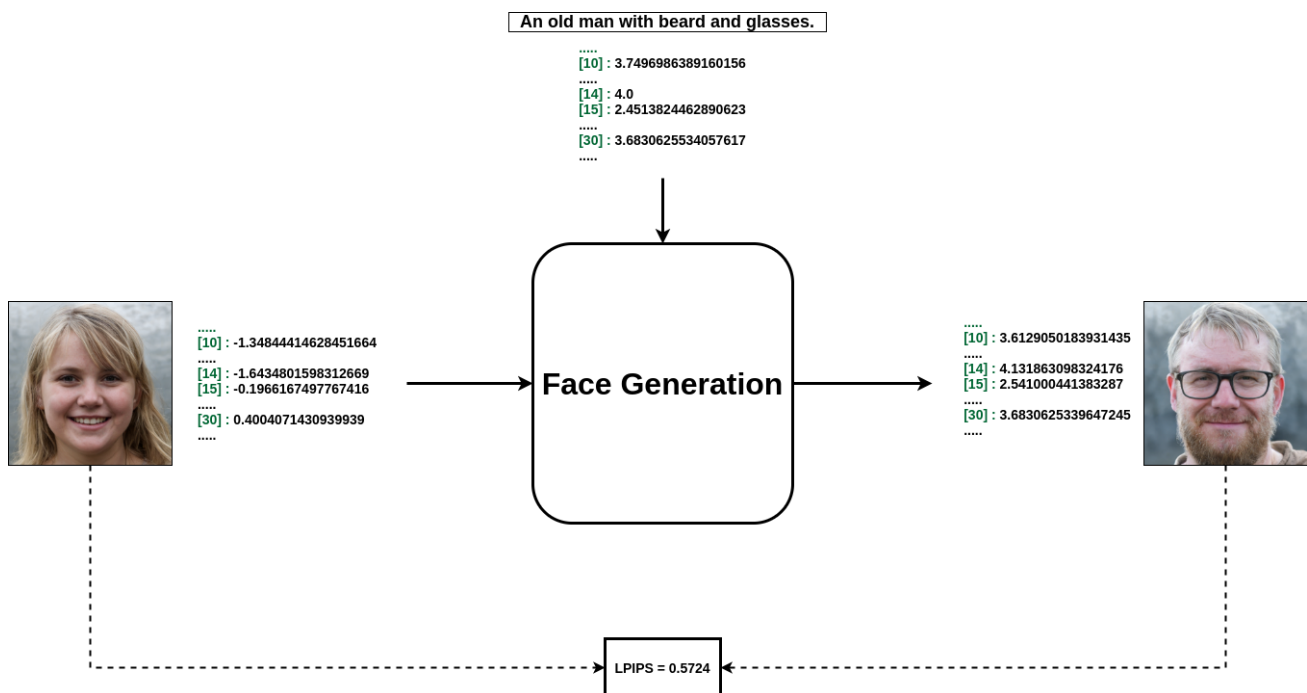


Figure 5.2: An example of the consistency in reaching the required facial attributes starting from initial random vector.

Input Text	Number of Navigated Directions	LPIPS
Female with chubby face	2	0.3905
Woman with long wavy hair	3	0.4412
Old black man with glasses	4	0.5023
Woman with lipstick and rosy cheeks	3	0.5107
Young black man with long hair and beard	5	0.5733
Old man with chubby face and glasses	4	0.4905

Table 5.2: LPIPS values against the number of navigated directions for sample text (lower is better).

5.2.1.5 Face Refinement

Since, we adapt StyleGAN2 for the refinement process, so this module testing is almost the same as **code-to-face generation** module. However, we focus more on the visuals of the *sequential navigation*, along with the new features related to the face morphology. Figure 5.3 shows the results of sequential navigation given an original synthetic face image. We tried to maintain the independence of sequential direction navigation as much as possible to keep the results visually acceptable. For example, in the first row, we convert from a *“young man with hear and no beard”* into an *“old bald man with beard”* by sequential navigation using *age*, *beard* and *baldness* directions. This yields much better visual results than many recent *attribute-editing GANs*.

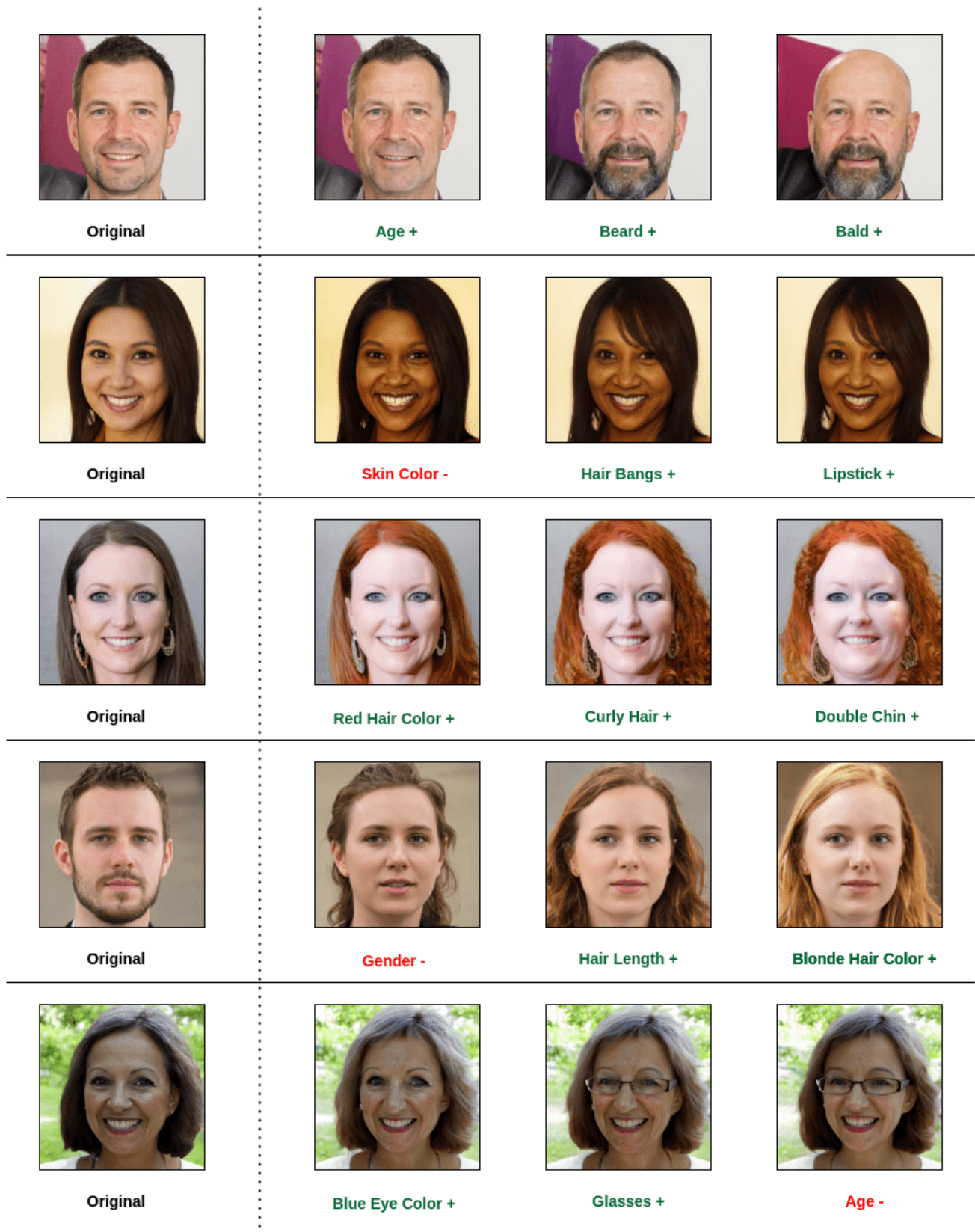


Figure 5.3: The results of sequential navigation along certain feature directions.

5.2.1.6 Multiple Head Poses Generation

5.2.1.7 Web Application

5.2.2 Integration Testing

The integration of the whole system into a web application is tested qualitatively and quantitatively. In this section, we show the visual results of our system in an *end-to-end* manner. However, we show the quantitative metrics in the comparison with previous work.

Figures 5.4 and 5.5 show correct results of **face generation** in an *end-to-end* manner using our final *web application*, which contains our complete pipeline. We can see that the results samples are of high accuracy and fidelity. However, of course, the system is *not* 100% accurate. From our extensive testing to the system, we notice *4types* of failures, described visually in figure 5.6. These failures are listed as follows :

1. Failures due to **contradicting facial features**. As in the *top-left* image, as *rosy cheeks* are very hard to capture with *black skin color*.
2. Failures due to **excessive navigation on directions**. As in the *top-right* image, which is very unclear and of low quality with many visual artifacts.
3. Failures due to **random initialization**. As in the *bottom-left* image, where sometimes bad initial latent vector can cause the output to be visually inconsistent with many visual artifacts.
4. Failures due to **sequential latent manipulation**. As in the *bottom-right* image, where navigation on one direction (*wavy hair*) cancel the navigation on another direction (*short hair*).



Young bald man with beard and glasses.



Young girl with blonde hair and blue eyes.



Young woman with brown wavy hair. She is putting on lipstick and wearing glasses.



Old black man with grey hair.

Figure 5.4: Samples of correctly generated face portrait from textual description.



A bald Asian man with beard. He is wearing glasses.



A white woman with short brown hair. She has hair bangs and is wearing glasses.



An old chubby man with double chin and rosy cheeks. He is wearing glasses and has receding hairline.



A young woman with blonde hair and rosy cheeks. She is putting on slight makeup and her hair is short.

Figure 5.5: Samples of correctly generated face portrait from textual description.



A black man with rosy cheeks and blonde hair.



A young woman with red hair and rosy cheeks. She is putting on lipstick.



A young man with long hair and beard.



A young woman with short wavy hair.

Figure 5.6: Samples of incorrectly generated face portrait from textual description.

5.3 Testing Schedule

Our testing process is scheduled as follows :

- First, we conducted the initial testing on the 3 core modules, while being iteratively designed, until they converged to decent results.
- We, then, did the integration testing on these modules separately.
- After so, the rest of the system modules were designed and tested separately.
- Finally, the whole system was integrated into a single web application and complete testing of the whole system functionalities was conducted.

5.4 Comparative Results to Previous Work

As mentioned before, we compare our results quantitatively with 3 baselines, using **FID score**, **average LPIPS** and **execution time**.

Table 5.3 and plot 5.7 show the comparison of our system to StyleGAN2 and Image2StyleGAN. We couldn't include Faces à la Carte in this comparison, as the authors didn't include it in the paper. Also, we couldn't replicate their work, as they didn't provide any specific details

about the implementation. We can see that as the number of test images increases, the overall quality of the images increases (*FID* score decreases). Our system performs better than Image2StyleGAN, but worse than the original StyleGAN2, because it is just image generation from random vector with no *latent manipulation* (fewer artifacts).

Test Size	StyleGAN2	Image2StyleGAN	Our System
50	129.15	179.58	151.58
100	104.29	150.54	132.54
200	95.25	136.45	114.45
300	92.44	134.04	113.04
400	86.37	121.59	100.59
500	80.09	119.02	99.02

Table 5.3: FID scores comparison on different number of test images (lower is better).



Figure 5.7: Plot of FID score of different pipelines against different number of test images (lower is better).

Next, we compare our system with Faces à la Carte using the only metric, the authors provided, which *LPIPS*. We can see from table 5.4 that our system yields lower overall *LPIPS*, which is better. However, we have higher error margin than Faces à la Carte.

Our System	Faces à la Carte
0.595±0.008	0.634±0.005

Table 5.4: LPIPS comparison with Faces à la Carte (lower is better).

Finally, table 5.5 shows the execution time of different core stages of our system. We can see that the overall *text-to-face generation* process takes only about 0.2 seconds. However, the web application can take from 1 second up to multiple seconds to generate a face portrait (from any description) depending on the *connectivity* with the server. Meanwhile, our *multiple head poses generation* module takes around 5 seconds.

Text Processing	Latent Manipulation	Face Generation
0.048	0.11	0.024

Table 5.5: The execution time of different stages of the core of our system (measured in seconds).

6 Conclusions and Future Work

6.1 Faced Challenges

6.2 Gained Experience

6.3 Conclusions

6.4 Future Work

References

- [1] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan, 2020.
- [2] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks, 2019.
- [3] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data, 2020.
- [4] Tianren Wang, Teng Zhang, and Brian Lovell. Faces à la carte: Text-to-face generation via attribute disentanglement, 2020.
- [5] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space?, 2019.
- [6] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018.
- [7] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions, 2014.
- [8] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric, 2018.

A Development Platforms and Tools

A.1 Hardware Platforms

A.2 Software Tools

B Use Cases

C User Guide

D Code Documentation

E Feasibility Study