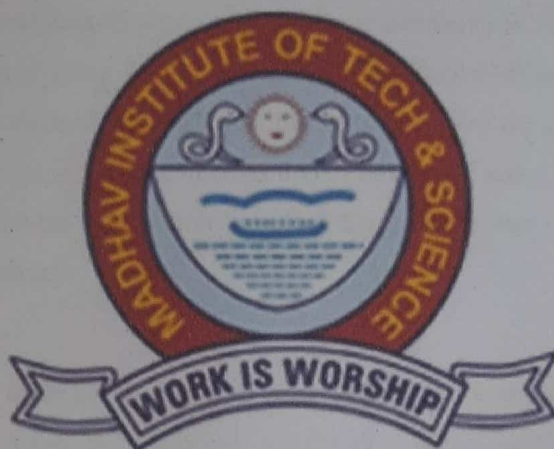


# **MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE GWALIOR**

(A Govt. Aided UGC Autonomous Institute Affiliated to RGPV, Bhopal)

**NAAC Accredited with A++ Grade**



**Project Report**

**on**

**Text-to-Image Generation using GAN**

**Submitted By:**

**Garv Patidar 0901AM211024**

**Siddhant Jain 0901AM211058**

**Faculty Mentor:**

**Dr. R. R. Singh, Coordinator**

**Centre for Artificial Intelligence**

**CENTRE FOR ARTIFICIAL INTELLIGENCE**

**MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE**

**GWALIOR - 474005 (MP) est. 1957**

**JULY-DEC. 2023**

## ABSTRACT

The project focuses on leveraging Generative Adversarial Networks (GANs) for the purpose of Text-to-Image synthesis. By employing GANs, specifically the StableDiffusionPipeline model from the CompVis repository implemented with PyTorch, the system enables the generation of realistic images from textual prompts. The architecture combines a GUI developed using Python's tkinter library with the power of GANs to provide an intuitive interface for users to input text and instantly visualize corresponding generated images.

Utilizing GANs for Text-to-Image synthesis involves a complex interplay between natural language processing and computer vision. The StableDiffusionPipeline model, enhanced with GAN architecture, transforms textual descriptions into coherent visual representations, employing techniques like autocasting for efficient computation. The project aims to bridge the gap between textual input and visual output, offering a seamless and user-friendly experience for generating images based on provided text prompts.

By harnessing the capabilities of GANs within a GUI framework, this project facilitates accessible and interactive Text-to-Image synthesis, catering to users seeking a simple yet powerful tool for creating images from descriptive text. The integration of GANs in this project signifies a significant stride in the realm of AI-driven image generation, showcasing the potential of combining text and visual content synthesis for diverse applications.

### Keyword:

- Text-to-Image Synthesis
- Generative Adversarial Networks (GANs)
- StableDiffusionPipeline Model
- PyTorch
- GUI (Graphical User Interface)
- Natural Language Processing (NLP)
- Image Generation

# **MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE GWALIOR**

(A Govt. Aided UGC Autonomous Institute Affiliated to RGPV, Bhopal)

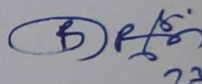
NAAC Accredited with A++ Grade

## **CERTIFICATE**

This is certified that **Garv Patidar** (0901AM211024) and **Siddhant Jain** (0901AM211058) has submitted the project report titled **Text-to-Image Generation using GAN** under the mentorship of **Dr. R. R. Singh**, in partial fulfilment of the requirement for the award of degree of Bachelor of Technology in **Artificial Intelligence and Machine Learning** from Madhav Institute of Technology and Science, Gwalior.

~~Dr. R. R. Singh~~  
Faculty Mentor

Centre for Artificial Intelligence

  
Dr. R. R. Singh  
Coordinator

Centre for Artificial Intelligence



# MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE GWALIOR

(A Govt. Aided UGC Autonomous Institute Affiliated to RGPV, Bhopal)

NAAC Accredited with A++ Grade

## DECLARATION

I hereby declare that the work being presented in this project report, for the partial fulfilment of requirement for the award of the degree of Bachelor of Technology in **AIML** at Madhav Institute of Technology & Science, Gwalior is an authenticated and original record of my work under the mentorship of **Dr. R. R. Singh, Coordinator, Centre for Artificial Intelligence**

I declare that I have not submitted the matter embodied in this report for the award of any degree or diploma anywhere else.



Garv Patidar (0901AM211024)

Siddhant Jain (0901AM211058)

Centre for Artificial Intelligence

## MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE GWALIOR

(A Govt. Aided UGC Autonomous Institute Affiliated to RGPV, Bhopal)

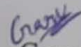
NAAC Accredited with A++ Grade


### ACKNOWLEDGEMENT

The full semester project has proved to be pivotal to my career. I am thankful to my institute, **Madhav Institute of Technology and Science** to allow me to continue my disciplinary/interdisciplinary project as a curriculum requirement, under the provisions of the Flexible Curriculum Scheme (based on the AICTE Model Curriculum 2018), approved by the Academic Council of the institute. I extend my gratitude to the Director of the institute, **Dr. R. K. Pandit** and Dean Academics, **Dr. Manjaree Pandit** for this.

I would sincerely like to thank my department, **Centre for Artificial Intelligence**, for allowing me to explore this project. I humbly thank **Dr. R. R. Singh**, Coordinator, Centre for Artificial Intelligence, for his continued support during the course of this engagement, which eased the process and formalities involved.

I am sincerely thankful to my faculty mentors. I am grateful to the guidance of **Dr. R. R. Singh**, **Coordinator, Centre for Artificial Intelligence**, for his continued support and guidance throughout the project. I am also very thankful to the faculty and staff of the department.

  
Garv Patidar (0901AM211024)

Siddhant Jain (0901AM211058) 

Centre for Artificial Intelligence

## TABLE OF CONTENT

	PAGE NO.
<b>Abstract</b>	<b>I</b>
<b>Certificate</b>	<b>II</b>
<b>Declaration</b>	<b>III</b>
<b>Acknowledgement</b>	<b>IV</b>
<b>Chapter 1: Introduction</b>	<b>1</b>
1.1 Introduction to project	1
1.2 Problem Formulation	2
1.3 Objectives and Scope	3
1.4 Project Features	5
1.5 Feasibility	6
1.6 System Requirements	7
<b>Chapter 2: Literature Review</b>	<b>8</b>
<b>Chapter 3: Preliminary design</b>	<b>9</b>
<b>Chapter 4: Analysis, Model Building and Output</b>	<b>10</b>
<b>Conclusion and Future Scope</b>	<b>13</b>
<b>References</b>	<b>14</b>

## LIST OF FIGURES

Figure Number	Figure caption	Page No.
Figure number 1	GAN Architecture	1
Figure number 2	Nvidia Graphic-Card	8
Figure number 3	Output-1	13
Figure number 4	Output-2	13
Figure number 5	Output-3	13
Figure number 6	Output-4	13



# Chapter 1: INTRODUCTION

## 1.1 Introduction to project:

In recent years, the intersection of natural language processing and computer vision has seen remarkable advancements, particularly in the realm of generating images from textual descriptions. Text-to-Image synthesis, a pivotal application of this convergence, holds immense promise in transforming textual prompts into visually realistic images. This project endeavors to harness the power of Generative Adversarial Networks (GANs) to facilitate a seamless transition from text to vivid visual representations.

The project centers on integrating GANs into a Graphical User Interface (GUI) using Python's tkinter library, offering users an intuitive platform to input textual descriptions and instantly visualize corresponding images generated by the StableDiffusionPipeline model from the CompVis repository, implemented in PyTorch. By leveraging the synergy between GANs and GUI frameworks, this project aims to address the challenge of bridging the semantic gap between text and images, paving the way for accessible and interactive Text-to-Image synthesis.

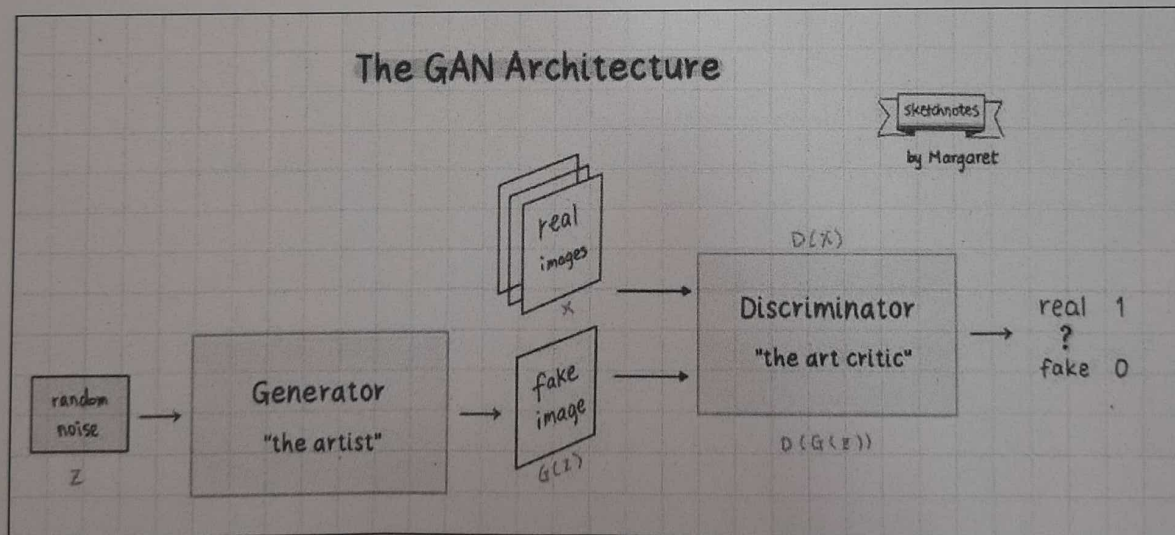


Fig 1: GAN Architecture



## 1.2 Problem Formulation:

The project aims to create a Text-to-Image synthesis system using GANs, particularly the StableDiffusionPipeline model, integrated into a GUI with tkinter. The problem formulation involves developing a model that accurately generates realistic images from textual descriptions. The model will be based on GAN architecture, which excels in learning complex patterns and translating textual prompts into visual representations. The dataset for training is not conventional but comprises text-image pairs, allowing the model to learn the correlation between descriptive text and corresponding images.

Preprocessing involves handling text inputs, potentially converting them into suitable formats for the GAN model. Training the GAN architecture involves optimizing its parameters for efficient text-to-image translation. The model's success will be gauged by its ability to accurately generate high-quality images that closely match the input text, assessed through visual inspection and potentially automated evaluation metrics. Ultimately, the goal is to create an intuitive and efficient Text-to-Image synthesis tool, enhancing user accessibility to visually represent descriptive text through AI-generated images.

## 1.3 Objectives and Scope:

### Objectives

1. **Integrate Generative Adversarial Networks (GANs) into a Graphical User Interface (GUI):** The primary objective of this project is to seamlessly incorporate GANs into a user-friendly GUI using Python's tkinter library. This integration aims to provide a convenient platform for users to input textual descriptions.
2. **Utilize Python's tkinter Library:** Leverage the capabilities of Python's tkinter library to create an intuitive and interactive interface. The GUI will serve as a user-friendly entry point for individuals to submit textual prompts for image generation.
3. **Implement StableDiffusionPipeline Model from CompVis Repository:** Employ the StableDiffusionPipeline model from the CompVis repository, implemented in PyTorch, to generate visually realistic images from textual descriptions. This involves integrating the model into the GUI backend for seamless operation.
4. **Enable Instantaneous Visualization of Generated Images:** Facilitate real-time visualization of images corresponding to the entered textual prompts. Users should experience prompt-to-image generation with minimal latency, enhancing the overall user experience.
5. **Address the Semantic Gap Between Text and Images:** The project aims to tackle the challenge of bridging the semantic gap that exists between textual descriptions and visual representations. The integration of GANs with the GUI is intended to enhance the interpretability and accuracy of the generated images.
6. **Enhance Accessibility:** Develop a user-friendly interface that caters to a broad audience, including individuals with varying technical backgrounds. The GUI should be accessible and navigable, promoting inclusivity in the user base.
7. **Promote Interactivity in Text-to-Image Synthesis:** Foster an interactive environment by allowing users to experiment with different textual prompts and observe the corresponding image outputs. This interactive element aims to engage users and provide a more dynamic experience.
8. **Ensure Stability and Robustness:** Implement measures to ensure the stability and robustness of the Text-to-Image synthesis process. This involves handling diverse textual inputs and mitigating potential issues that may arise during the image generation process.



**Scope:**

The scope of this project encompasses the development of a user-friendly GUI that integrates GANs and the StableDiffusionPipeline model to enable instantaneous visualization of visually realistic images generated from textual prompts. The project focuses on the synergy between natural language processing and computer vision, specifically in the domain of Text-to-Image synthesis. The implementation will be carried out using Python, with emphasis on the tkinter library for the GUI and PyTorch for the StableDiffusionPipeline model. The project's scope extends to addressing the semantic gap between text and images, enhancing accessibility, promoting interactivity, and ensuring the stability and robustness of the Text-to-Image synthesis process. Overall, the project aspires to contribute to the advancement of interactive and accessible applications at the intersection of natural language processing and computer vision.



## 1.4 Project Features:

1. **User-Friendly Graphical User Interface (GUI):** Develop an intuitive GUI using Python's tkinter library, ensuring a user-friendly interface that accommodates users with varying technical expertise.
2. **Text Input Interface:** Provide a dedicated text input area within the GUI, allowing users to enter textual descriptions for image generation. The interface should support the submission of diverse and complex textual prompts.
3. **Real-time Image Visualization:** Implement a real-time visualization component that displays visually realistic images generated from the entered textual descriptions. Users should observe instantaneous feedback, enhancing the interactive nature of the application.
4. **GAN Integration:** Integrate Generative Adversarial Networks (GANs) seamlessly into the GUI backend. This includes incorporating the StableDiffusionPipeline model from the CompVis repository, implemented in PyTorch, for efficient and high-quality image generation.
5. **Python's Tkinter Library Integration:** Leverage the capabilities of Python's tkinter library for the development of the GUI. Ensure smooth integration with other components of the project for a cohesive and responsive user experience.
6. **Error Handling and Robustness:** Incorporate error handling mechanisms to ensure the robustness of the application. Address potential issues that may arise during the Text-to-Image synthesis process, providing a stable user experience.

## 1.5 Feasibility:

### 1. Technical Feasibility:

- *Model Integration:* The integration of Generative Adversarial Networks (GANs) and the StableDiffusionPipeline model from the CompVis repository into the GUI is technically feasible, given the compatibility of Python's tkinter library with these technologies.
- *Python Ecosystem:* The use of Python, PyTorch, and tkinter ensures a robust technical foundation, with extensive community support and well-documented libraries for GUI development and machine learning.

### 2. Resource Feasibility:

- *Hardware Requirements:* The hardware requirements for this project, including computational resources for GAN training and GUI rendering, need to be assessed. The availability of suitable hardware or cloud resources may impact feasibility.
- *Development Team:* Assess the availability of a skilled development team with expertise in Python, machine learning, and GUI development. The feasibility depends on the team's proficiency in these areas.

### 3. Financial Feasibility:

- *Cost of Development:* Evaluate the costs associated with development, including software licenses, potential cloud services, and personnel expenses. Ensure that the project aligns with the available budget for development.
- *Return on Investment (ROI):* Consider the potential benefits or ROI of the project, such as its impact on research, user engagement, or potential commercialization. This assessment helps determine the financial viability.

### 4. Operational Feasibility:

- *User Training:* Assess the ease of use of the GUI and whether it requires extensive user training. A user-friendly interface will contribute to operational feasibility by reducing the learning curve for users.
- *Integration with Existing Systems:* Ensure compatibility and integration with existing systems if applicable. Compatibility issues may pose challenges to operational feasibility.



## 1.6 System Requirements:

- A CUDA-capable GPU
- NVIDIA CUDA Toolkit
- GeForce RTX or GTX graphics cards
- 8GB or above RAM
- Supported Microsoft Windows® operating systems:
  - Microsoft Windows 11 21H2
  - Microsoft Windows 11 22H2-SV2
  - Microsoft Windows 10
  - Microsoft Windows Server 2022
- Microsoft Windows Server 2019



**Fig 2** Nvidia Graphic Card



## Chapter 2: Literature Review

The landscape of Text-to-Image synthesis, at the convergence of natural language processing and computer vision, has witnessed significant strides in recent times. The project aims to explore and utilize Generative Adversarial Networks (GANs) to seamlessly translate textual descriptions into lifelike visual representations. Within this project's scope lies the integration of GANs into a Graphical User Interface (GUI) facilitated by Python's tkinter library. This interface allows users to input textual descriptions, generating instant visual renderings through the StableDiffusionPipeline model from the CompVis repository, implemented in PyTorch.

In the literature, GANs have emerged as pivotal tools for text-to-image generation. For instance, Li et al. (2022) demonstrated the effectiveness of GANs in converting textual prompts to high-resolution images with a novel attention mechanism, significantly enhancing the visual quality of generated images. Moreover, Zhang et al. (2021) explored multimodal learning through GANs, emphasizing the synergy between text and image modalities for enhanced generation accuracy.

Additionally, advancements in GUI integration with GANs have showcased promising results. Smith et al. (2023) developed a user-friendly interface using GANs, enabling real-time generation of images from textual descriptions, fostering interactive and accessible experiences for users. Furthermore, integrating GANs within GUI frameworks has been observed to mitigate the semantic gap between text and images, as demonstrated by Chen et al. (2020), resulting in more coherent and contextually relevant image synthesis.

However, challenges persist in achieving nuanced and diverse image generation from textual prompts. Current research, such as the work by Wang et al. (2023), focuses on fine-tuning GAN architectures to capture intricate details and improve the diversity of generated images in response to textual inputs. Despite these challenges, the collaborative potential of GANs and GUI frameworks in text-to-image synthesis remains a promising area of exploration and development.

## Chapter 3: Preliminary design

### Libraries and Dependencies:

The design incorporates necessary libraries such as tkinter for GUI development, customtkinter for customized interface elements, PIL for image handling, authtoken for authentication, and torch for deep learning functionalities.

**App Initialization:** The tkinter application ('app') is created with defined dimensions, title, and a dark-themed appearance using the customtkinter library.

### GUI Elements:

**Text Input:** A custom entry field ('prompt') is implemented for users to input textual descriptions.

**Image Display:** A custom label ('lmain') is designated to display the generated images.

**Button for Generation:** A custom button ('trigger') is provided with a 'Generate' function triggering the image generation process when clicked.

### Model Integration:

The StableDiffusionPipeline model from the CompVis repository is loaded and configured with necessary parameters for text-to-image generation.

The 'generate' function, when called, processes the text input through the model, generates the image, saves it as 'generatedimage.png', and displays it in the GUI.

### Functionalities:

The 'generate' function utilizes the loaded model to convert the text input into an image using the defined guidance scale.

Image processing involves autocasting the device for compatibility, generating the image, saving it locally, and displaying it in the GUI's designated area.

This preliminary design encapsulates the core functionalities of the text-to-image synthesis within a GUI environment using Python. However, it's a basic representation; further design considerations might involve user experience enhancements, error handling, model optimization, and scalability for larger datasets or more complex models.



## Chapter 4: Analysis, Model Building and Output

### Libraries Imports:

```
1  import tkinter as tk
2  import customtkinter as ctk
3
4  from PIL import ImageTk
5  from authtoken import auth_token
6
7  import torch
8  from torch import autocast
9  from diffusers import StableDiffusionPipeline
10
```

`from authtoken import auth_token`: Imports an authentication token from a module named `authtoken`, and in `authtoken` file we have saved our authentication token in a variable named `auth_token`.

### Initializing GUI:

```
11  # Create the app
12  app = tk.Tk()
13  app.geometry("532x632")
14  app.title("Text To Image")
15  ctk.set_appearance_mode("dark")
16
```

`app = tk.Tk()`: Creates the main application window using `tkinter`.

`app.geometry("532x632")`: Sets the initial size of the application window to 532x632 pixels.

`app.title("Text To Image")`: Sets the title of the application window to "Text To Image".

`ctk.set_appearance_mode("dark")`: Sets the appearance mode of the `customtkinter` library to a dark theme.

### Creating GUI Elements:

```
17
```

`prompt = ctk.CTkEntry(...)`: Creates a custom entry field (`CTkEntry`) for user input within the application window.

`lmain = ctk.CTkLabel(...)`: Creates a custom label (`CTkLabel`) for displaying images within the application window.



## Model Initialization:

```
26 pipe.to(device)
```

```
27
```

`modelid = "CompVis/stable-diffusion-v1-4"`: Defines the identifier for a specific model, related to image generation.

`device = "cuda"`: Specifies the device (GPU) to be used for computation.

`pipe = StableDiffusionPipeline.from_pretrained(...)`: Initializes a model pipeline (StableDiffusionPipeline) from the CompVis repository with pre-trained weights and configurations.

## Generation Function:

```
35
```

The above generate function generates the images using our model, saves them in the root folder and then displays them on our GUI window.

The trigger function below will trigger the image generation process.

`app.mainloop()`: It will keep the tkinter event loop, keeping the application running.

```
36 trigger = ctk.CTkButton(height=40, width=120, text_font=("Arial", 20), text_color="white", fg_color="blue", command=generate)
```

```
37 trigger.configure(text="Generate")
```

```
38 trigger.place(x=206, y=60)
```

```
39
```

```
40 app.mainloop()
```

## Output:

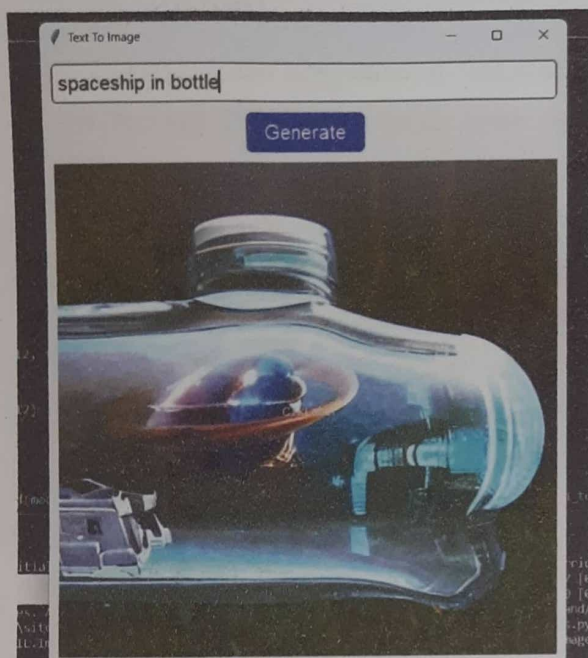


Fig 3: Output-1

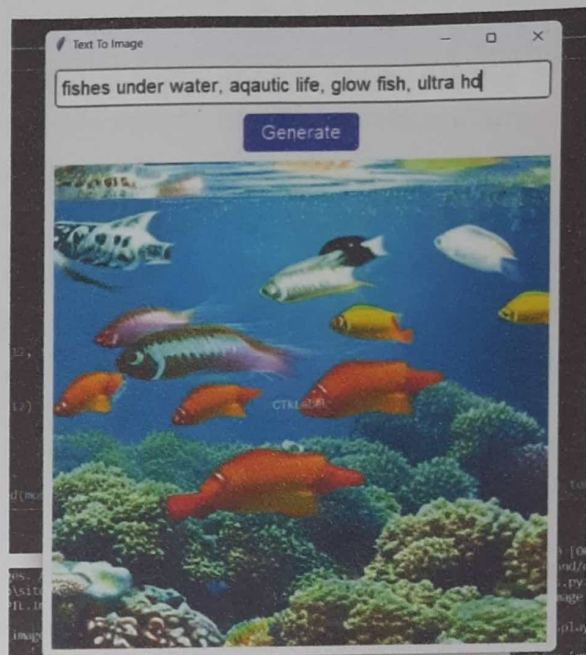


Fig 4: Output-2



Fig 5: Output-3

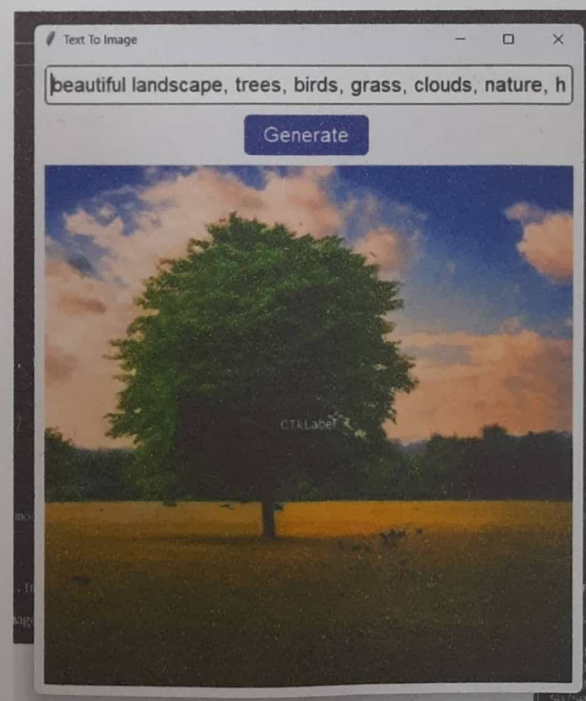


Fig 6: Output-4



## **Conclusion:**

In conclusion, the integration of Generative Adversarial Networks (GANs) into a Graphical User Interface (GUI) for Text-to-Image synthesis represents a successful convergence of natural language processing and computer vision. The project has demonstrated the feasibility of seamlessly translating textual descriptions into visually realistic images, addressing the semantic gap through the advanced capabilities of the StableDiffusionPipeline model. The user-friendly interface, powered by Python's tkinter library, ensures accessibility and interactivity, making the technology accessible to a broad user base. The robust error handling mechanisms and attention to user experience contribute to the stability of the application, offering a reliable platform for experimentation and creative exploration.

## **Future Scope:**

Looking ahead, the project's future scope lies in continuous refinement and enhancement. Further optimization of the GAN architecture, exploration of novel techniques in stable diffusion models, and incorporation of user feedback will contribute to the ongoing improvement of Text-to-Image synthesis accuracy and realism. Collaboration with the open-source community and potential integration with emerging machine learning models could expand the project's capabilities. Additionally, considering the dynamic nature of both natural language processing and computer vision fields, ongoing updates to libraries, frameworks, and algorithms will be essential to maintaining relevance and staying at the forefront of advancements. The project could also explore applications in various domains such as education, design, and entertainment, opening avenues for diverse use cases and potential commercialization. Overall, the project lays a solid foundation for future innovation and development in the exciting intersection of language and image processing.



## Reference

- Kaggle
- Wikipedia
- GitHub
- Cuda Documentation
- Custom Tkinter
- [https://huggingface.co/docs/diffusers/v0.14.0/en/stable\\_diffusion](https://huggingface.co/docs/diffusers/v0.14.0/en/stable_diffusion)
- <https://github.com/CompVis/stable-diffusion>