

# Text to Image Generator Using Python

Dr. Ramesh Kumar C<sup>1</sup>, Fahad<sup>2</sup>, Abrar<sup>3</sup>, Rehan<sup>4</sup>

<sup>1</sup>Associate Professor, Department of Computer Science & Engineering, Deccan College of Engineering and Technology, Hyderabad, India

<sup>234</sup>UG Students, Department of Computer Science & Engineering, Deccan College of Engineering and Technology, Hyderabad, India

**Abstract**—This project presents the design and development of a fully offline-capable Text-to-Image Generator that utilizes advanced generative AI techniques. By integrating Stable Diffusion XL (SDXL), Gradio, Python, and Docker, this system enables users to convert natural language prompts into high-quality images without relying on internet connectivity. This solution addresses key limitations of existing cloud-based tools such as privacy concerns, dependency on online infrastructure, and recurring costs. The system supports multiple image generation styles (realistic and anime), image download options, history tracking, and performance metrics. It serves as an efficient and scalable tool for educational, creative, and secure environments where offline usage is essential.

**Index Terms**—Offline AI, Stable Diffusion XL, Text-to-Image, Gradio, Image Generation, Generative AI, Docker, Safetensors

## 1. INTRODUCTION

Text-to-image synthesis is a rapidly evolving area in generative AI, where systems translate natural language descriptions into coherent and realistic images. Existing solutions such as DALL·E, MidJourney, and Stable Diffusion often depend on internet connectivity, cloud computing, and external servers. These dependencies pose challenges in terms of accessibility, data privacy, and operational cost.

To address these concerns, this project proposes a fully offline Text-to-Image Generator using Stable Diffusion XL, Gradio UI, and Docker containerization. Users can operate the application entirely on local machines while benefiting from robust image quality, prompt flexibility, and UI intuitiveness. The system supports both real-time generation and style-switching modes (anime and realistic), providing a customizable and powerful tool for users with limited or no internet access.

### 1.2 Problem Statement:

Current text-to-image generation platforms are heavily dependent on internet connectivity and cloud-based APIs, which limits their usability in offline or secure environments. These systems often pose challenges related to data privacy, latency, recurring service costs, and dependency on external infrastructure. In educational institutions, research labs, or remote areas, such limitations hinder creative exploration and practical deployment. Additionally, most existing solutions lack customization options and control over system resources. This project aims to address these issues by offering a fully offline, secure, and user-friendly text-to-image generation tool.

### 1.3 OBJECTIVE

The primary objective of this project is to develop a fully offline Text-to-Image Generator that allows users to convert natural language prompts into high-quality images without relying on internet connectivity. By integrating Stable Diffusion XL with a locally hosted interface using Gradio, the system aims to provide a seamless and intuitive experience for both technical and non-technical users. The project emphasizes privacy, speed, and reliability, making it ideal for environments where data security and offline functionality are critical. Additionally, it seeks to support multiple image styles, enhance usability through customization features, and offer performance metrics to monitor generation efficiency.

### 1.4 Scope

The scope of this project includes the design, development, and deployment of an offline-capable Text-to-Image Generator that operates entirely without internet access after installation. It covers

the integration of Stable Diffusion XL for generating high-quality images, a user-friendly interface built with Gradio, and secure model handling using Safetensors. The system supports multiple features such as real-time and anime-style image generation, image history tracking, download options, and basic performance monitoring. This project is intended for users in educational, research, or secure environments who need reliable AI tools without cloud dependency. Future expandability, including new model integrations and UI enhancements, is also considered within the system's architectural framework.

## 2. LITERATURE SURVEY

The design and development of an offline text-to-image generator is grounded in a diverse range of research areas including generative modeling, edge computing, interface design, and privacy-enhanced deployment. The following subsections review the key technological foundations and highlight the research gaps this project addresses.

### 2.1 Evolution of Text-to-Image Synthesis

Early advancements in text-to-image generation were led by GAN-based models such as StackGAN and AttnGAN. These systems aimed to generate images based on textual input but struggled with issues like low visual fidelity, limited contextual understanding, and unstable training dynamics. Later, transformer-based architectures like OpenAI's DALL·E and Tsinghua's CogView introduced autoregressive modeling and multi-head attention, significantly improving coherence between input text and output image. These models laid the groundwork for modern systems by enabling better semantic representation and image realism.

### 2.2 Diffusion Models

Diffusion-based generative models emerged as a robust alternative to GANs, particularly with the introduction of Denoising Diffusion Probabilistic Models (DDPM) by Ho et al. These models reverse a noise-adding process to generate images from random noise, offering better training stability and high-quality outputs. Stable Diffusion further refined this approach by introducing latent diffusion techniques, optimizing performance and computational efficiency. The release of Stable Diffusion XL (SDXL) enhanced capabilities with support for high-resolution generation and improved

interpretation of complex prompts, making it ideal for offline, high-fidelity synthesis.

### 2.3 Edge Deployment and Offline AI

Recent literature in edge computing and private AI emphasizes the value of offline inference systems to ensure user privacy, reduce latency, and maintain control over computational resources. Unlike cloud-based services, local models remove dependency on internet connectivity and prevent third-party data access. Tools like Docker are widely adopted for deploying AI models locally in isolated, reproducible environments. Their platform-independent nature makes them especially useful for consistent and secure offline deployments.

### 2.4 Secure Model Loading with Safetensors

Security in model loading is critical for offline systems. Traditional .ckpt files can pose security risks by enabling the execution of arbitrary Python code during deserialization. Safetensors, developed by Hugging Face, addresses this by storing only raw tensor data, eliminating execution risk. This format supports faster, deterministic, and secure model loading, making it suitable for environments where user safety and reproducibility are priorities.

### 2.5 User Interface with Gradio

User experience is central to AI accessibility. Gradio has become a popular framework for building ML interfaces due to its intuitive Python API, support for offline usage, and minimal front-end development requirements. It allows developers to create interactive UIs that work locally without needing port exposure or online dependencies. For projects like this one, Gradio's ability to quickly bind models to browser-accessible apps makes it an optimal choice for building usable, responsive interfaces for non-technical users.

### 2.6 Research Gaps and Contribution

While the field has progressed rapidly, a significant gap remains in developing fully offline, secure, and user-friendly systems for text-to-image generation. Most existing research and tools focus on cloud-based implementations with limited attention to customization, security, and offline functionality. Additionally, features like generation history, resolution customization, and user-controlled style

switching (e.g., anime vs. realistic) are seldom integrated into a single package. This project addresses these gaps by delivering a complete offline solution tailored for privacy, performance, and usability.

### 3. EXISTING SYSTEM

Existing text-to-image generation systems, such as OpenAI's DALL·E, MidJourney, and Google's Imagen, rely heavily on internet access and cloud-based infrastructure. These platforms offer high-quality image outputs by using advanced transformer or diffusion-based models. Users submit text prompts via web interfaces, and the image generation occurs on remote servers. While powerful, these systems require continuous internet connectivity and often store user data for analysis or improvement of models. Additionally, customization options and access to underlying model architecture are usually restricted in these commercial services.

#### 3.1 Limitations of Existing System

- Constant Internet Dependency

Most current text-to-image generation platforms require a stable internet connection, making them unusable in remote, secure, or offline environments such as defense sectors, rural areas, or isolated research labs.

- Data Privacy Risks

Online platforms may collect and log user prompts and generated images. This can lead to unintentional exposure of sensitive or personal data, especially when the system is used in enterprise or academic settings.

- Limited Customization Options

Users have minimal control over the model's behavior, generation parameters, or output style. Advanced configurations, such as adjusting image resolution or choosing between realistic and anime outputs, are often unavailable.

- Closed-Source or Proprietary Models

Many platforms, such as MidJourney or DALL·E, operate on proprietary architectures. Users cannot inspect or modify the code or understand the inner workings of the model.

- Subscription and Usage Costs

High-quality image generation often comes at a cost, including subscription fees, credit-based usage models,

or premium memberships, which can be restrictive for students, researchers, or small organizations.

- Lack of Offline Functionality

These platforms are hosted on remote servers, and thus cannot be used without internet access. This makes them unsuitable for environments with strict security protocols or limited connectivity.

- Scalability Issues in Free Tiers

Free or trial versions of online services usually impose restrictions on resolution, number of generations, or processing priority, limiting usability for heavy or large-scale projects.

- Limited User Interface Control

Users often cannot alter the layout, behavior, or functionalities of the platform interface, preventing integration with other tools or personalized workflows.

- Inability to Track or Manage History

Online systems generally do not offer features like image generation history, prompt tracking, or version control. This makes it hard for users to revisit or manage past outputs effectively.

- Vulnerability to Server Downtime or API Changes

Reliance on external services makes users dependent on the availability and stability of cloud providers. Any downtime, model change, or API deprecation can halt or break existing workflows.

#### 3.2 PROPOSED SYSTEM

The proposed system is a fully offline, secure, and customizable Text-to-Image Generator that leverages Stable Diffusion XL for high-quality image generation based on natural language prompts. Unlike existing online platforms, this system runs entirely on local hardware, ensuring data privacy, operational control, and zero dependency on internet connectivity.

At its core, the system utilizes latent diffusion models to efficiently generate images while maintaining superior visual quality. The inclusion of the Safetensors format guarantees secure and deterministic model loading without any risk of

hidden code execution, making it ideal for sensitive or restricted environments.

To enhance user experience, the system integrates a Gradio-based local interface, allowing users to interact with the model through a clean, browser-accessible UI.

Key functionalities include:

- Multiple image generation per prompt (1, 2, or 4 options)
- Switchable styles like Realistic and Anime Download options (single, bulk, with metadata)
- Viewing generation history and performance metrics (time, count, etc.)
- No background data collection or online access

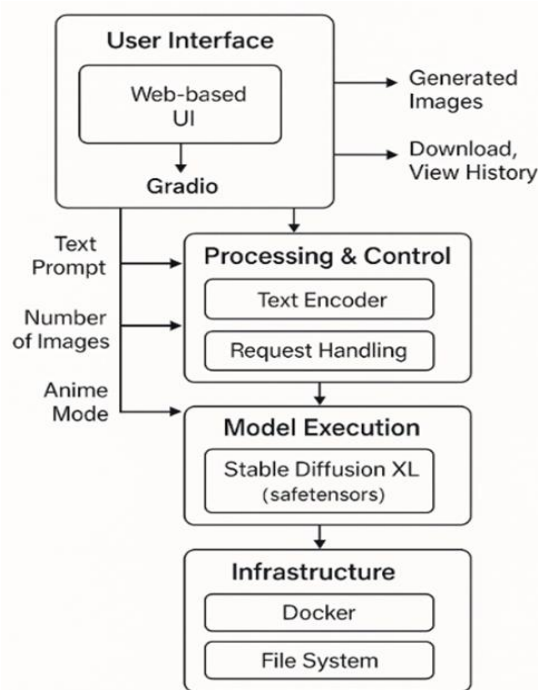


Fig 1 Architecture

The deployment is containerized using Docker, ensuring easy installation and consistency across various operating systems. This also allows smooth portability and scalability in offline environments.

### 3.2.1 Advantages of the Proposed System

- **Fully Offline Operation:** The system runs without any internet connection, ensuring complete data privacy and making it suitable for secure or remote environments.

- **Enhanced Security with Safetensors:** By using the Safetensors format, the model loads safely without executing arbitrary code, protecting against potential security risks associated with traditional model files.

- **High-Quality and Customizable Outputs:** Powered by Stable Diffusion XL, the system generates high-resolution, realistic or anime-style images with improved prompt understanding and visual coherence.

- **User-Friendly Interface via Gradio:** A simple, browser-based interface allows both technical and non-technical users to generate images, view history, and download results with ease.

- **Platform Independence and Easy Deployment:** Docker-based deployment ensures consistent performance across different systems and simplifies installation without complex setup requirements.

### 3.4 TOOLS AND TECHNOLOGIES USED

The development of the offline Text-to-Image Generator integrates a combination of modern tools and technologies across machine learning, deployment, and user interface domains:

#### 1. Stable Diffusion XL (SDXL)

A state-of-the-art text-to-image model that uses latent diffusion for efficient and high-quality image generation. It supports detailed prompts, higher resolution, and better visual realism.

#### 2. Safetensors

A secure and fast model serialization format developed by Hugging Face. It ensures safe model loading by preventing execution of embedded code, making it ideal for offline and private deployments.

#### 3. Gradio

An open-source Python library used to create a local, interactive, browser-based user interface. It enables real-time input/output handling with minimal coding, perfect for non-technical users.

#### 4. Docker

A containerization platform used to package the application with all dependencies, ensuring consistent behavior across different operating systems and simplifying deployment.

### 5. Python

The core programming language used to build and integrate model execution, UI components, logic control, and image handling.

### 6. JavaScript

Used to enhance the frontend interactivity of the Gradio interface, including features like real-time updates, history display, and download options.

### 7. HTML & CSS

Used within Gradio templates for customizing UI layout and styling, ensuring a user-friendly and clean interface design.

### 8. NumPy & PIL

Python libraries used for image manipulation, processing, and array-based operations required during generation and download phases.

## 3.5 ALGORITHMS USED

The primary algorithm used in this project is the Latent Diffusion Model (LDM), a powerful and efficient form of generative modeling. It is based on the Denoising Diffusion Probabilistic Model (DDPM) but optimized for performance by operating in a lower-dimensional latent space.

**1. Denoising Diffusion Probabilistic Model (DDPM)**  
Originally proposed by Ho et al. (2020), DDPM is a generative model that creates data (such as images) by reversing a gradual noising process. The model is trained to predict and remove noise from an image over many steps, effectively learning how to “denoise” a random noise vector into a coherent image.

**Forward Process:** Adds small amounts of Gaussian noise to an image over  $T$  time steps until it becomes pure noise.

**Reverse Process:** Learns to reconstruct the image by removing noise step by step, generating a new image from randomness.

### 2. Latent Diffusion Model (LDM)

Instead of applying DDPM directly to high-resolution pixel images (which is computationally expensive), LDM first compresses images into a lower-dimensional latent space using an autoencoder. Diffusion is then performed in this latent space, significantly reducing memory and processing requirements.

**Text Conditioning:** A CLIP-based encoder translates user prompts into embeddings that guide the diffusion process.

**Latent Space Processing:** The model denoises within a compressed space, then decodes the final result back into a high-quality image.

**Efficiency & Speed:** This method drastically improves speed and allows for high-resolution image generation even on limited hardware.

### 3. Classifier-Free Guidance

To improve alignment between generated images and the input text, the model uses classifier-free guidance. It generates two predictions — one conditioned on the text prompt and one without — and combines them to produce images that better match the user’s intent.

## 4. CONCLUSION AND FUTURE ENHANCEMENT

### Conclusion

The development of the Text-to-Image Generator Without Internet demonstrates the feasibility of deploying advanced AI models locally without compromising functionality or user experience. By integrating Stable Diffusion XL with technologies like Python, Gradio, Docker, and Safetensor, the system enables offline generation of high-quality, contextually relevant images from natural language prompts. The user-friendly interface offers features such as image style toggling, generation history, and download options, making the application accessible for a wide range of users.

This offline solution addresses critical limitations found in cloud-based platforms, including internet dependency, data privacy concerns, and ongoing service costs. Its modular and scalable architecture allows for future enhancements and adaptation to various hardware environments. The project thus provides a robust, secure, and flexible tool suitable for educational, research, and creative professionals needing offline generative capabilities.

### Future Enhancements

Several avenues exist to further improve the system. Expanding model support by including multiple Stable Diffusion versions and domain-specific fine-

tuned models would increase versatility. Enhancements to prompt input, such as suggestion features and filtered prompt history, can improve user creativity. Advanced controls like resolution adjustment, negative prompts, and seed locking would refine output customization. Incorporating basic image editing and local inpainting features could allow post-generation modifications. Introducing multi-user management with role-based access and private histories would support collaborative environments. Performance monitoring tools and intelligent resource management can optimize system efficiency. Additionally, an optional internet-connected mode for updates and community features, as well as offline voice-to-text prompt input, would enhance usability and future-proof the application.

[20] <https://www.intel.com/content/www/us/en/artificial-intelligence/posts/edge-ai.html>

## REFERENCES

- [1] Ho, J., Jain, A., & Abbeel, P. (2020). Denoising Diffusion Probabilistic Models.
- [2] <https://arxiv.org/abs/2006.11239>
- [3] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-Resolution Image Synthesis with Latent Diffusion Models.
- [4] <https://arxiv.org/abs/2112.10752>
- [5] Stability AI. (2023). Stable Diffusion XL (SDXL).
- [6] <https://stability.ai/news/stable-diffusion-xl-release>
- [7] Hugging Face. (2023). Safetensors Documentation.
- [8] <https://huggingface.co/docs/safetensors>
- [9] Abid, A., et al. (2019). Gradio: Hassle-Free Sharing and Testing of ML Models in the Wild.
- [10] <https://gradio.app>
- [11] Merkel, D. (2014). Docker: Lightweight Linux Containers for Consistent Development and Deployment.
- [12] <https://www.docker.com/resources/what-container>
- [13] Zhang, H., et al. (2017). StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks.
- [14] <https://arxiv.org/abs/1612.03242>
- [15] Xu, T., et al. (2018). AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks.
- [16] <https://arxiv.org/abs/1711.10485>
- [17] Ramesh, A., et al. (2021). Zero-Shot Text-to-Image Generation (DALL·E).
- [18] <https://arxiv.org/abs/2102.12092>
- [19] Edge AI Deployment – Intel AI Blog.