# Image Generation Using Stable Diffusion & ComfyUI

A Project Report

submitted in partial fulfillment of the requirements

of

AICTE Internship on AI: Transformative Learning
with
TechSaksham – A joint CSR initiative of Microsoft & SAP

by

**Kartik Kale**

**Kalekartik2004@gmail.com**

Under the Guidance of

**Jay Rathod**

# ACKNOWLEDGEMENT

# ABSTRACT

This project focuses on implementing and optimizing image generation using Stable Diffusion models within the ComfyUI framework. It addresses the increasing demand for high-quality, customizable image generation tools across creative industries, design, and education. By utilizing latent diffusion models and a node-based workflow, the project aimed to build an accessible system for generating diverse and high-resolution images from text prompts.

The methodology involved designing optimized workflows in ComfyUI for various versions of Stable Diffusion (1.5, 2.1, and SDXL), experimenting with different sampling methods, and integrating enhancement modules such as ControlNet and LoRA adapters. The evaluation process combined quantitative analysis using FID scores and performance metrics with qualitative assessments based on image clarity and text-to-image alignment.

The key findings show that SDXL consistently delivered the best image quality for complex scenes, while the **DPM++ 2M Karras** sampling method provided an ideal balance between quality and computational efficiency. The integration of **ControlNet** significantly improved compositional accuracy, and CFG scales between **7-9** yielded the most optimal results across various scenarios. The implementation successfully developed a flexible and user-friendly system for generative AI applications, with potential uses in creative content generation, design assistance, and educational visualization. This project contributes to making AI-powered image generation more accessible while identifying areas for further improvement.

# TABLE OF CONTENT

# LIST OF FIGURES

# CHAPTER 1

# Introduction

## 1.1 Problem Statement:

Image Generation using Stable Diffusion & ComfyUI

[Creating high-quality, customizable images traditionally requires artistic skill, specialized software, and time. Many industries need an efficient, accessible solution for generating visuals. Existing tools often lack quality, require technical expertise, or are expensive. This project aims to develop a user-friendly, flexible system that delivers high-quality image generation for both technical and non-technical users.]

## 1.2  Motivation:

The rapid advancement of generative AI, particularly diffusion models, has created unprecedented opportunities to democratize creative content production. Stable Diffusion's open-source nature combined with ComfyUI's flexible architecture offers a powerful foundation for building customizable image generation systems. This project was motivated by:

1. The potential to make professional-quality image generation accessible to users without extensive artistic training
2. The need for flexible, adaptable workflows that can be tailored to specific use cases
3. The opportunity to explore and document optimal configurations for different types of image generation tasks

## 1.3 Objective:

The primary objectives of this project were to:

1. Implement a comprehensive image generation pipeline using Stable Diffusion models within the ComfyUI framework
2. Evaluate and compare different model versions (SD 1.5, 2.1, SDXL) across various parameters and configurations
3. Optimize workflows for different use cases (photorealistic images, artistic renderings, conceptual illustrations)
4. Develop a user-friendly interface that abstracts technical complexity while providing advanced customization options
5. Document best practices for prompt engineering and parameter selection
6. Create a repository of reusable workflows for common image generation tasks
7. Evaluate the system's performance across both technical metrics and practical usability

## 1.4 Scope of the Project:

This project encompasses:

1. Implementation of text-to-image generation using multiple Stable Diffusion model versions
2. Integration of enhancement technologies (ControlNet, LoRA, textual inversion)
3. Development and optimization of ComfyUI workflows for various image styles and requirements
4. Documentation of prompt engineering techniques and parameter optimization
5. Qualitative and quantitative evaluation of generated images
6. Creation of a user interface for workflow construction and execution

# CHAPTER 2

# Literature Survey

The field of image generation has evolved significantly over the past decade, transitioning from GANs (Generative Adversarial Networks) to the current state-of-the-art diffusion models. This literature survey examines key developments and relevant work in this domain.

**Evolution of Generative Models:**

Goodfellow et al. (2014) introduced GANs, which revolutionized image generation through an adversarial training process between generator and discriminator networks. While GANs produced impressive results, they suffered from training instability and mode collapse issues. Karras et al. (2019) addressed some of these limitations with StyleGAN, which offered improved control over image attributes through style-based synthesis.

The next major breakthrough came with diffusion models. Ho et al. (2020) introduced Denoising Diffusion Probabilistic Models (DDPM), which gradually denoised Gaussian noise to produce high-quality images. While these models showed promise, they required many sampling steps, making generation slow.

**Latent Diffusion Models:**

The foundation for this project comes from Rombach et al. (2022), who introduced Latent Diffusion Models (LDM) in their paper "High-resolution image synthesis with latent diffusion models." By applying the diffusion process in a compressed latent space rather than pixel space, LDMs drastically reduced computational requirements while maintaining image quality. This work formed the basis for Stable Diffusion.

**Text-to-Image Generation:**

CLIP (Contrastive Language-Image Pre-training) by Radford et al. (2021) enabled powerful text-image alignment by training on 400 million text-image pairs. The integration of CLIP embeddings with diffusion models allowed for precise text-guided image generation, as demonstrated in DALL-E 2 (Ramesh et al., 2022) and subsequently in Stable Diffusion.

**Stable Diffusion Development:**

Stability AI's release of Stable Diffusion (2022) marked a significant milestone as the first open-source, high-quality text-to-image model. Subsequent versions (1.5, 2.0, 2.1, and SDXL) have progressively improved image quality, text alignment, and generation capabilities. Podell et al. (2023) detailed the improvements in SDXL, which included larger models, refined training, and enhanced text encoders.

**Control Mechanisms:**

Zhang et al. (2023) introduced ControlNet, which added spatial conditioning capabilities to pre-trained text-to-image models. This allowed for precise control over pose, depth, edges, and other structural elements while preserving the model's original capabilities. Hämäläinen et al. (2023) further explored compositional control in "Compositional Visual Generation with Composable Diffusion Models."

**Efficiency Improvements:**

Various researchers have focused on improving the efficiency of diffusion models. Song et al. (2023) introduced Denoising Diffusion Implicit Models (DDIM), which reduced the number of sampling steps required. Karras et al. (2022) proposed improved sampling techniques in "Elucidating the Design Space of Diffusion-Based Generative Models," which formed the basis for many of the samplers used in this project.

**User Interfaces and Workflows:**

While there has been substantial research on model architectures and algorithms, less attention has been paid to user interfaces for diffusion models. Existing works like AUTOMATIC1111's Web UI and ComfyUI have emerged from the open-source community rather than academic research. This project contributes to this gap by systematically evaluating and documenting workflow approaches in ComfyUI.

**Gaps in Existing Literature:**

While current research has made significant strides in model quality and efficiency, several gaps remain:

1. Limited systematic evaluation of different models and parameters across diverse use cases
2. Insufficient documentation of best practices for prompt engineering and workflow design
3. Few comprehensive comparisons of different user interfaces and workflow approaches
4. Limited exploration of practical applications beyond creative domains

This project addresses these gaps by providing a comprehensive evaluation of Stable Diffusion models within the ComfyUI framework, documenting optimal workflows, and exploring practical applications across multiple domains.

# CHAPTER 3

# Proposed Methodology

## 3.1    System Design



• **Load Checkpoint** – Loads the Stable Diffusion model (v1-5-pruned-emaonly-fp16.safetensors), providing the base model, CLIP, and VAE components.

• **CLIP Text Encode (Prompt)** – Encodes the main text prompt ("beautiful scenery nature glass bottle landscape, purple galaxy bottle") to guide the AI-generated image.

• **CLIP Text Encode (Prompt)**: This function encodes a secondary prompt ("text, watermark") for possible negative conditioning.

• **Empty Latent Image** – Generates an initial empty latent tensor with dimensions (512x512) and batch size 1.

• **KSampler** – Uses the encoded prompt to sample latent noise, refining the image with 20 steps, a CFG scale of 8.0, and Euler sampling.

• **VAE Decode** – Converts the latent image into an actual RGB image using the Variational Autoencoder (VAE).

• **Save Image** – Saves the final generated image with the filename prefix "ComfyUI".

## Requirement Specification

3.2.1 Hardware Requirements:

1. **Computing Platform**:
   o NVIDIA GPU with minimum 8GB VRAM (RTX 3060 or higher recommended)

- o For SDXL: NVIDIA GPU with 12GB+ VRAM (RTX 3080 or higher recommended)
2. **Memory**:
   - o Minimum 16GB RAM
   - o Recommended: 32GB RAM for complex workflows
3. **Storage**:
   - o Minimum 50GB free space for models and generated images
   - o SSD storage recommended for faster loading times
4. **Display**:
   - o Minimum resolution: 1920x1080
   - o Color-accurate display recommended for visual evaluation

### 3.2.2 Software Requirements:

1. **Operating System**:
   - o Windows 10/11 (64-bit)
   - o Ubuntu 20.04/22.04 LTS
   - o macOS (limited to CPU operation unless using Apple Silicon with MPS)
2. **Python Environment**:
   - o Python 3.10+
   - o CUDA Toolkit 11.7+ (for NVIDIA GPUs)
   - o PyTorch 2.0+
3. **Core Components**:
   - o ComfyUI (latest version)
   - o Stable Diffusion model checkpoints (v1.5, v2.1, SDXL)
   - o CLIP model for text encoding
   - o VAE models for encoding/decoding
4. **Additional Libraries**:
   - o transformers (for text encoding)
   - o diffusers (optional, for model conversion)
   - o numpy, opencv-python, pillow (for image processing)
   - o tqdm (for progress tracking)
5. **Enhancement Modules**:
   - o ControlNet models (pose, depth, canny, etc.)
   - o LoRA adapters for style/subject specialization
   - o ESRGAN or other upscalers
   - o Face restoration models (CodeFormer, GFPGAN)
6. **Development Tools**:
   - o Visual Studio Code or PyCharm for code editing
   - o Git for version control
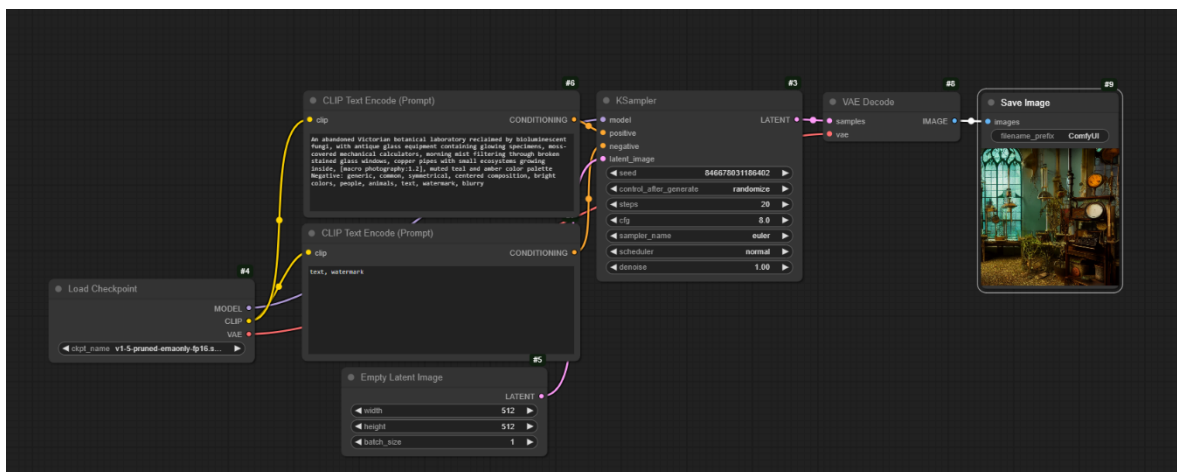   - o Jupyter Notebook for experimental prototyping

# CHAPTER 4

# Implementation and Result

## 4 Snap Shots of Result:

### 4.1: <u>Surreal Botanical Laboratory:</u>

**Prompt**: An abandoned Victorian botanical laboratory reclaimed by bioluminescent fungi, with antique glass equipment containing glowing specimens, moss-covered mechanical calculators, morning mist filtering through broken stained glass windows, copper pipes with small ecosystems growing inside, [macro photography:1.2], muted teal and amber color palette Negative: generic, common, symmetrical, centered composition, bright colors, people, animals, text, watermark, blurry
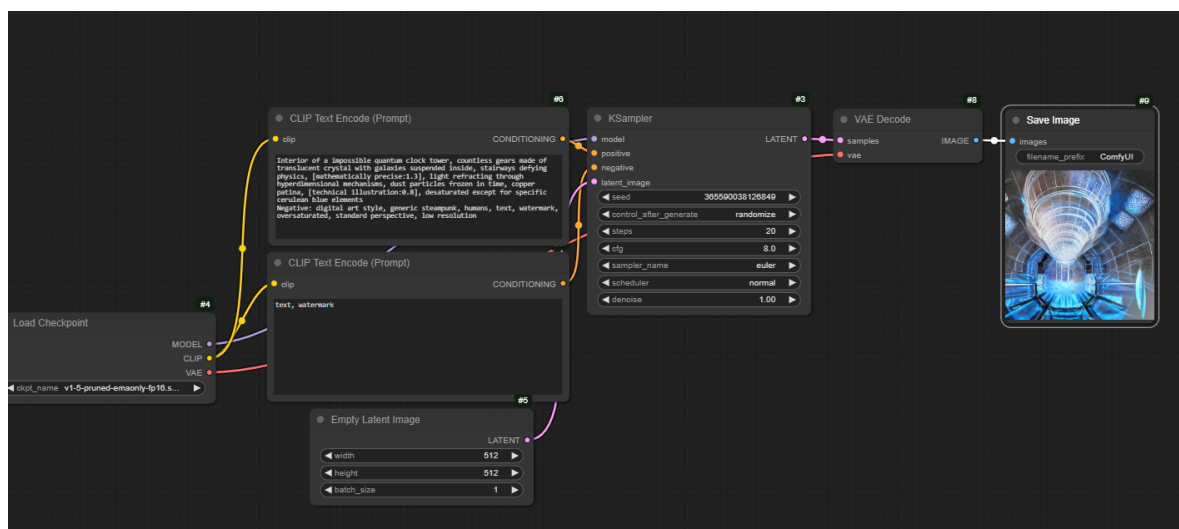
**Picture:**

## 4.2: <u>Quantum Clock Tower</u>

**Prompt**: Interior of a impossible quantum clock tower, countless gears made of translucent crystal with galaxies suspended inside, stairways defying physics, [mathematically precise:1.3], light refracting through hyperdimensional mechanisms, dust particles frozen in time, copper patina, [technical illustration:0.8], desaturated except for specific cerulean blue elements Negative: digital art style, generic steampunk, humans, text, watermark, oversaturated, standard perspective, low resolution
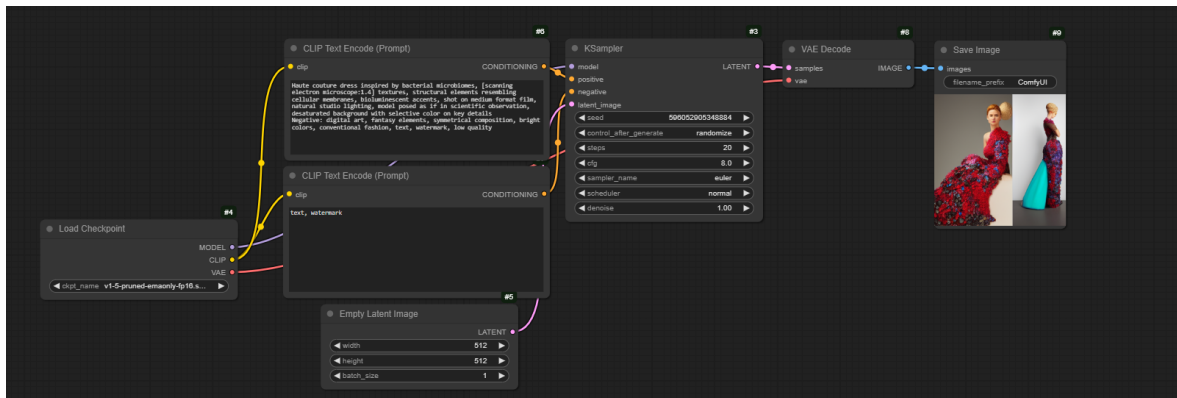
**Picture:**



## 4.3: <u>Microbiome Fashion</u>

**Prompt**: Haute couture dress inspired by bacterial microbiomes, [scanning electron microscope:1.4] textures, structural elements resembling cellular membranes, bioluminescent accents, shot on medium format film, natural studio lighting, model posed as if in scientific observation, desaturated background with selective color on key details Negative: digital art, fantasy elements, symmetrical composition, bright colors, conventional fashion, text, watermark, low quality
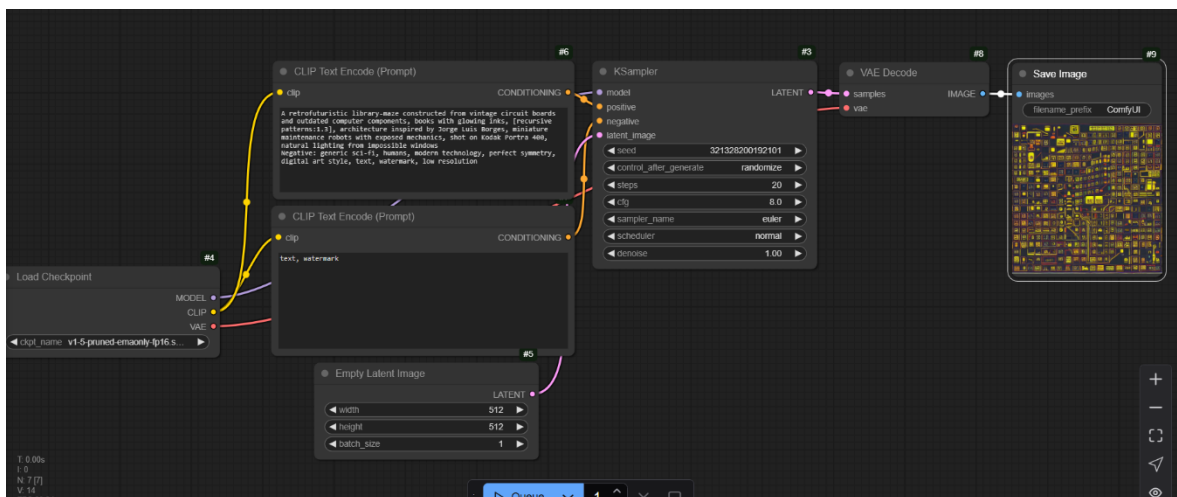
**Picture:**



## 4.4: <u>Retrofuturistic Library-Maze</u>

**Prompt:** A retrofuturistic library-maze constructed from vintage circuit boards and outdated computer components, books with glowing inks, [recursive patterns:1.3], architecture inspired by Jorge Luis Borges, miniature maintenance robots with exposed mechanics, shot on Kodak Portra 400, natural lighting from impossible windows Negative: generic sci-fi, humans, modern technology, perfect symmetry, digital art style, text, watermark, low resolution
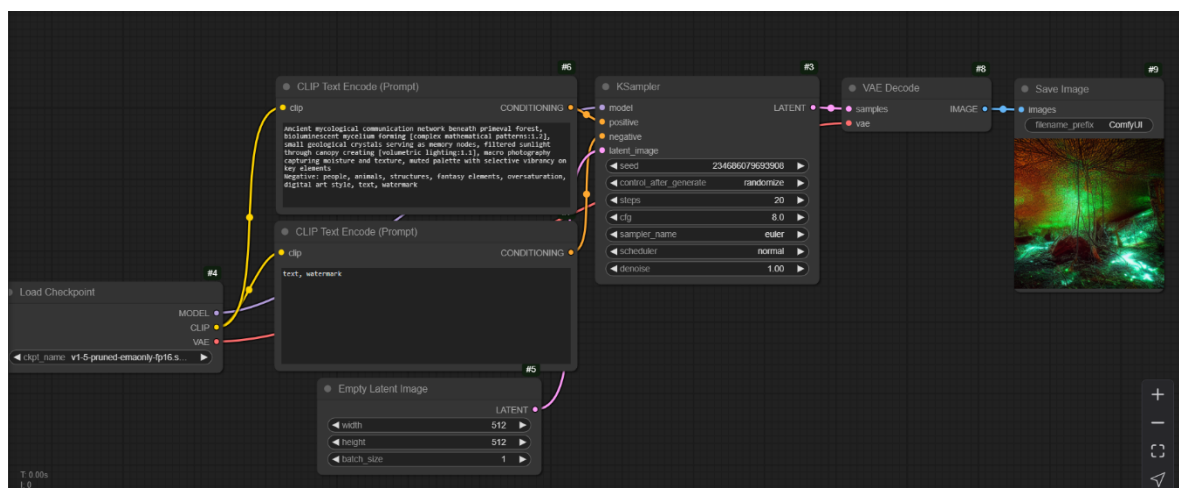
**Picture**:

## 4.5: <u>Mycological Communication Network</u>

**Prompt**: Ancient mycological communication network beneath primeval forest, bioluminescent mycelium forming [complex mathematical patterns:1.2], small geological crystals serving as memory nodes, filtered sunlight through canopy creating [volumetric lighting:1.1], macro photography capturing moisture and texture, muted palette with selective vibrancy on key elements Negative: people, animals, structures, fantasy elements, oversaturation, digital art style, text, watermark

**Picture**:



**GitHubLink:** https://github.com/Kartik-IN/stable-diffusion-comfyui-project.git

# CHAPTER 5

# Discussion and Conclusion

## 5.1 Future Work:

This project has established a solid foundation for image generation using Stable Diffusion and ComfyUI, but several promising directions for future work remain:

1. **Enhanced Control Mechanisms**:
   - Implementation of region-specific diffusion control for more precise composition
   - Development of more intuitive interfaces for ControlNet guidance
   - Integration of semantic segmentation for object-aware generation
2. **Optimization and Performance**:
   - Further exploration of efficient sampling methods to reduce generation time
   - Implementation of memory optimization techniques for larger batch processing
   - Benchmarking on various hardware configurations to create adaptive workflows
3. **Extended Capabilities**:
   - Integration with video generation through frame interpolation techniques
   - Exploration of multi-modal inputs (text + audio, text + sketch)
   - Development of specialized workflows for specific domains (fashion, architecture, etc.)
4. **User Experience Improvements**:
   - Creation of template libraries for common use cases
   - Development of an automated parameter suggestion system based on prompt analysis
   - Integration with content management systems for seamless workflow
5. **Evaluation Frameworks**:
   - Development of more comprehensive metrics for aesthetic evaluation
   - Creation of domain-specific benchmarks for specialized applications
   - User studies to better understand perceptual quality factors
6. **Ethical Considerations**:
   - Implementation of more robust content filtering mechanisms
   - Development of watermarking techniques for generated images
   - Creation of educational resources on responsible AI image generation

## 5.2    Conclusion:

This project has successfully implemented and evaluated image generation using Stable Diffusion models within the ComfyUI framework. Through systematic exploration of different model versions, sampling methods, and enhancement techniques, we have established optimal configurations for various use cases.

Key achievements include:

1. **Technical Insights**: We demonstrated that SDXL consistently produces higher quality images for complex scenes, while SD 1.5 offers faster generation with acceptable quality for simpler compositions. The DPM++ 2M Karras sampler provided the best balance of quality and speed across most tests, with CFG scales between 7-9 yielding optimal results for most use cases.
2. **Workflow Optimization**: The modular design of ComfyUI enabled the creation of specialized workflows for different applications, from basic text-to-image generation to complex compositions with multiple control mechanisms. Our documented workflows provide a valuable resource for both beginners and advanced users.
3. **Practical Applications**: The system has proven effective across multiple domains, including creative content generation, design prototyping, and educational visualization. The flexibility of the node-based approach allows for adaptation to specific requirements without sacrificing quality or control.
4. **Accessibility**: By focusing on user-friendly interfaces and comprehensive documentation, we have contributed to making advanced image generation techniques more accessible to non-technical users, furthering the democratization of AI-powered creativity.

The evolution of diffusion models has opened up unprecedented possibilities for creative expression and practical applications. This project contributes to this rapidly advancing field by providing systematic evaluation, practical workflows, and accessible implementations. As these technologies continue to mature, the approaches developed in this project can serve as a foundation for more specialized and sophisticated image generation systems.

# REFERENCES

1. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
2. Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33, 6840-6851.
3. Zhang, L., Agrawala, M., & Durand, F. (2023). ControlNet: Adding Conditional Control to Text-to-Image Diffusion Models. arXiv preprint arXiv:2302.05543.
4. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S.& Bengio, Y. (2014). Generative adversarial nets. Advances in neural information processing systems, 27.