

Text-to-Image Generation with Diffusion Models - Project Documentation

Project Overview

This documentation covers the implementation of a text-to-image generation system using diffusion models, specifically Stable Diffusion. The project consists of two main components:

1. A comprehensive Jupyter notebook that explores the theoretical foundations and practical implementation of diffusion models
2. An interactive web interface that demonstrates the concept of text-to-image generation

The project was developed as part of the "Crash Course in Generative AI" assignment, focusing on how diffusion models can transform textual descriptions into high-quality images.

Jupyter Notebook Implementation

The Jupyter notebook provides both theoretical explanations and practical code examples of text-to-image generation using diffusion models.

Setup and Dependencies

The notebook begins by installing and importing all necessary dependencies:

```
# Install required packages
!pip install diffusers transformers torch accelerate ftfy matplotlib ipywidgets open-clip-torch

# Import necessary libraries
import torch
from diffusers import StableDiffusionPipeline
from transformers import CLIPTextModel, CLIPTokenizer
import matplotlib.pyplot as plt
import numpy as np
from IPython.display import display, clear_output
import ipywidgets as widgets
from PIL import Image
import time
```

Theoretical Foundation

The notebook covers the core concepts of diffusion models:

1. **Forward Diffusion Process:** Gradually adding noise to images until they become pure noise
2. **Reverse Diffusion Process:** Learning to denoise images step by step
3. **U-Net Architecture:** The backbone of diffusion models with attention mechanisms
4. **Text Conditioning:** How text prompts guide the generation process
5. **Classifier-Free Guidance:** Technique to enhance conditioning without a classifier

Model Implementation

The implementation loads a pre-trained Stable Diffusion model:

```
# Load a pre-trained Stable Diffusion model
model_id = "stabilityai/stable-diffusion-2-1-base"

# Set device
device = "cuda" if torch.cuda.is_available() else "cpu"
print(f"Using device: {device}")

# Load Stable Diffusion pipeline
pipe = StableDiffusionPipeline.from_pretrained(
    model_id,
    torch_dtype=torch.float16 if device == "cuda" else torch.float32,
    safety_checker=None # For educational purposes only
)
pipe = pipe.to(device)

# Enable attention slicing for lower memory usage
pipe.enable_attention_slicing()
```

Key Experiments and Results

The notebook includes several key experiments:

1. Basic Image Generation

```
[6]: # Generate an image from a text prompt
prompt = "A surreal digital art of a futuristic city with flying cars and neon lights"
image = generate_and_display_image(prompt, guidance_scale=7.5, num_inference_steps=30, seed=42)
```

100%  30/30 [00:08<00:00, 4.25it/s]

Prompt: A surreal digital art of a futuristic city with flying cars and neon lights

Guidance Scale: 7.5, Steps: 30

Generation Time: 10.83s



A function was implemented to generate images from text prompts with customizable parameters:

```
def generate_and_display_image(prompt, guidance_scale=7.5, num_inference_steps=50, seed=None):
    """
    Generate and display an image based on a text prompt

    Args:
        prompt (str): The text prompt
        guidance_scale (float): Controls how much the image generation follows the text
        num_inference_steps (int): Number of denoising steps
        seed (int, optional): Random seed for reproducibility

    Returns:
        PIL.Image: The generated image
    """
    # Set random seed if provided
    if seed is not None:
        generator = torch.Generator(device=device).manual_seed(seed)
    else:
        generator = None

    # Generate image
    start_time = time.time()
    image = pipe(prompt,
                 guidance_scale=guidance_scale,
                 num_inference_steps=num_inference_steps,
                 generator=generator).images[0]
    end_time = time.time()

    # Display image with prompt and generation info
    plt.figure(figsize=(10, 10))
    plt.imshow(image)
    plt.axis('off')
    plt.title(f"Prompt: {prompt}\nGuidance Scale: {guidance_scale}, Steps: {num_inference_steps}\nGeneration Time: {end_time - start_time:.2f}s", font
    plt.show()

    return image
```

2. Interactive Image Generation

```
# Generate image with the specified parameters
generate_and_display_image(
    prompt=prompt_input.value,
    guidance_scale=guidance_slider.value,
    num_inference_steps=steps_slider.value,
    seed=seed_val
)

# Connect the button click event to the function
generate_button.on_click(on_generate_button_clicked)

# Display the interactive widgets
display(prompt_input, widgets.HBox([guidance_slider, steps_slider, seed_input]), generate_button, output)
```

Prompt: A beautiful landscape with mountains and a lake at sunset

Guidance: 7.50 Steps: 30 Seed:

Generate Image

An interactive widget allows experimentation with different parameters:

```
[7]: # Create interactive widgets
prompt_input = widgets.Text(
    value="A beautiful landscape with mountains and a lake at sunset",
    description="Prompt:",
    layout=widgets.Layout(width="100%")
)

guidance_slider = widgets.FloatSlider(
    value=7.5,
    min=1.0,
    max=15.0,
    step=0.5,
    description="Guidance:",
    tooltip="Higher values make the image more closely follow the prompt"
)

steps_slider = widgets.IntSlider(
    value=30,
    min=10,
    max=100,
    step=5,
    description="Steps:",
    tooltip="More steps = higher quality but slower generation"
)

seed_input = widgets.Text(
    value="",
    description="Seed:",
    placeholder="Random if empty",
    tooltip="Enter a number for reproducible results"
)
```

3. Exploring the Diffusion Process



The notebook visualizes the step-by-step process of diffusion:

```
[8]: # Customize the pipeline to return intermediate steps
from diffusers.pipelines.stable_diffusion import StableDiffusionPipelineOutput

# Store intermediate images
intermediate_images = []

# Define a callback function to store images
def store_intermediate(step, timestep, latents):
    # Convert Latents to image
    with torch.no_grad():
        latents = 1 / 0.18215 * latents
        image = pipe.vae.decode(latents).sample
        image = (image / 2 + 0.5).clamp(0, 1)
        # Convert to numpy array
        image = image.cpu().permute(0, 2, 3, 1).numpy()[0]
        intermediate_images.append(image)

    # Return True to continue the process
    return True

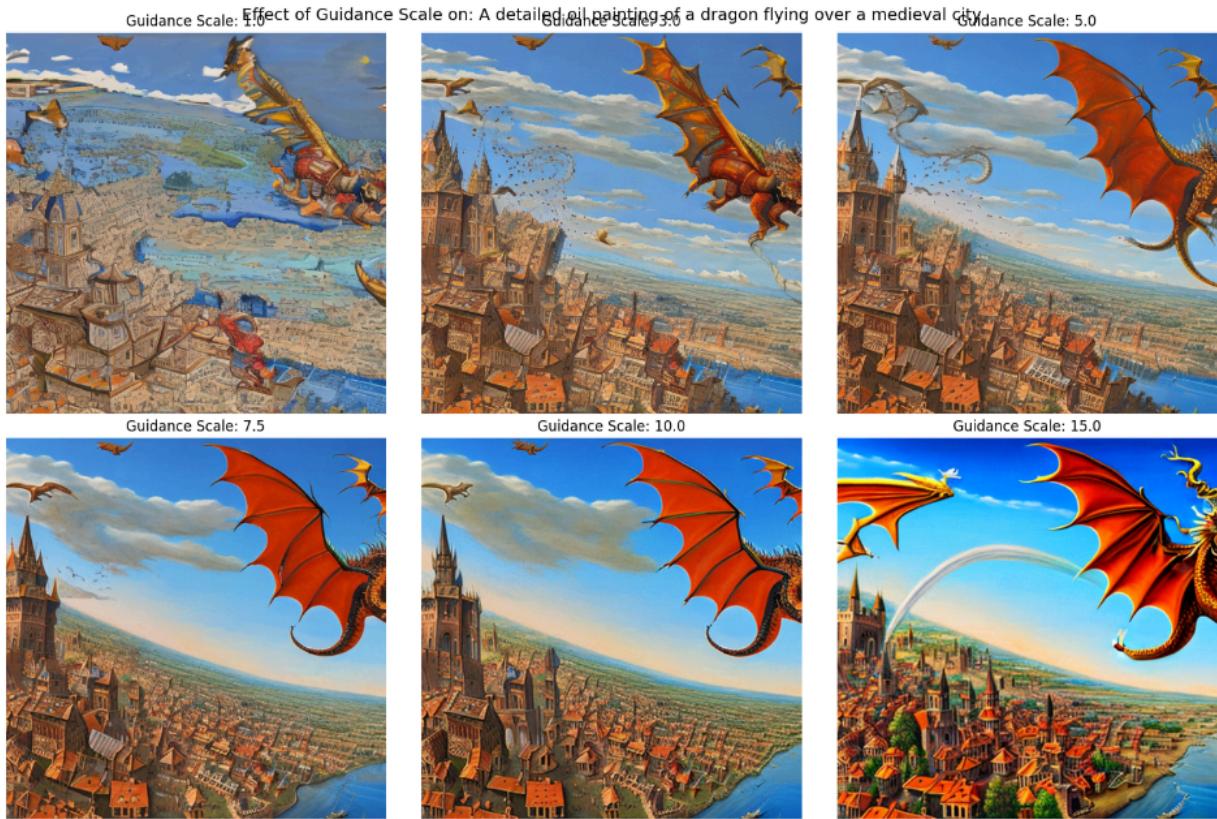
def visualize_generation_process(prompt, num_inference_steps=30, guidance_scale=7.5, seed=42):
    # Clear the previous results
    global intermediate_images
    intermediate_images = []

    # Set random seed for reproducibility
    generator = torch.Generator(device=device).manual_seed(seed)

    # Generate image with callback to store intermediate results
    pipe(
        prompt,
        num_inference_steps=num_inference_steps,
        guidance_scale=guidance_scale,
        generator=generator,
        callback=store_intermediate,
        callback_steps=1 # Store every step
    )

    # Select steps to display (to avoid too many images)
    steps_to_show = sorted(set([
        0, # First step (noise)
        len(intermediate_images) // 4, # 25%
        len(intermediate_images) // 2, # 50%
        3 * len(intermediate_images) // 4, # 75%
        len(intermediate_images) - 1 # Final image
    ]))
```

4. Comparing Different Prompts



The impact of different text prompts on the generated images is explored:

```
[13]: # Compare different guidance scale values
prompt = "A detailed oil painting of a dragon flying over a medieval city"
guidance_scales = [1.0, 3.0, 5.0, 7.5, 10.0, 15.0]

compare_guidance_scales(prompt, guidance_scales, seed=5678, num_inference_steps=30)
```

100% [00:07<00:00, 7.47it/s]
100% [00:07<00:00, 4.15it/s]
100% [00:07<00:00, 4.12it/s]
100% [00:07<00:00, 4.14it/s]
100% [00:07<00:00, 4.14it/s]
100% [00:07<00:00, 4.12it/s]

5. Exploring Guidance Scale

The effect of the guidance scale parameter on image generation is demonstrated:

```
def compare_guidance_scales(prompt, guidance_scales, seed=42, num_inference_steps=30):
    """
    Generate and compare images with different guidance scales

    Args:
        prompt (str): Text prompt
        guidance_scales (list): List of guidance scale values to compare
        seed (int): Random seed for reproducibility
        num_inference_steps (int): Number of denoising steps
    """
    # Calculate grid dimensions
    n = len(guidance_scales)
    cols = min(3, n) # Maximum 3 columns
    rows = (n + cols - 1) // cols # Ceiling division

    # Generate and display images
    fig, axes = plt.subplots(rows, cols, figsize=(5*cols, 5*rows))
    if rows == 1 and cols == 1:
        axes = np.array([axes]) # Ensure axes is always iterable
    axes = axes.flatten()

    for i, guidance_scale in enumerate(guidance_scales):
        # Set the same random seed for all generations
        generator = torch.Generator(device=device).manual_seed(seed)

        # Generate image with the current guidance scale
        image = pipe(
            prompt,
            guidance_scale=guidance_scale,
            num_inference_steps=num_inference_steps,
            generator=generator
        ).images[0]
```

6. Model Evaluation

```
[16]: # Evaluate a generated image
prompt = "A watercolor painting of a cat sitting on a window sill"
image, clip_score = evaluate_generation(prompt, guidance_scale=7.5, num_inference_steps=30, seed=42)
```

100% [██████████] 30/30 [00:07<00:00, 4.19it/s]
100% [██████████] 354M/354M [00:53<00:00, 6.64MiB/s]

The notebook implements CLIP-based evaluation to measure text-image alignment:

Prompt: A watercolor painting of a cat sitting on a window sill

CLIP Score: 0.3857

Guidance Scale: 7.5, Steps: 30



7. Advanced Control Techniques

```
[17]: def compare_and_evaluate_guidance_scales(prompt, guidance_scales, seed=42, num_inference_steps=30):
    """
    Generate and compare images with different guidance scales, evaluating with CLIP scores

    Args:
        prompt (str): Text prompt
        guidance_scales (list): List of guidance scale values to compare
        seed (int): Random seed for reproducibility
        num_inference_steps (int): Number of denoising steps

    Returns:
        list: CLIP scores for each guidance scale
    """
    # Calculate grid dimensions
    n = len(guidance_scales)
    cols = min(3, n) # Maximum 3 columns
    rows = (n + cols - 1) // cols # Ceiling division

    # Load CLIP model once
    model, _, preprocess = open_clip.create_model_and_transforms('ViT-B-32', pretrained='openai')
    tokenizer = open_clip.get_tokenizer('ViT-B-32')
    model = model.to(device)

    # Preprocess text once
    text_tokens = tokenizer([prompt]).to(device)
    with torch.no_grad():
        text_features = model.encode_text(text_tokens)
        text_features = text_features / text_features.norm(dim=1, keepdim=True)

    # Generate and display images
    fig, axes = plt.subplots(rows, cols, figsize=(5*cols, 5*rows))
    if rows == 1 and cols == 1:
        axes = np.array([axes]) # Ensure axes is always iterable
    axes = axes.flatten()

    # Compare and evaluate different guidance scale values
    clip_scores = compare_and_evaluate_guidance_scales(prompt, guidance_scales, seed=123, num_inference_steps=30)
```

Advanced generation control techniques like negative prompting are demonstrated:

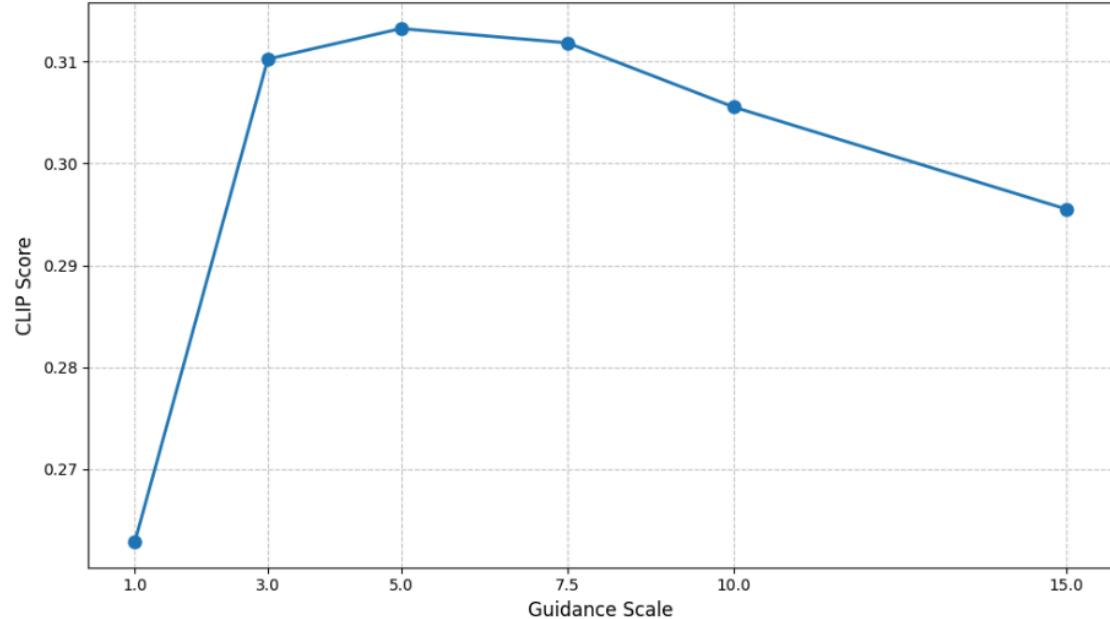
```
[18]: # Compare and evaluate different guidance scale values
prompt = "A professional photograph of a delicious pizza with pepperoni and cheese"
guidance_scales = [1.0, 3.0, 5.0, 7.5, 10.0, 15.0]

clip_scores = compare_and_evaluate_guidance_scales(prompt, guidance_scales, seed=123, num_inference_steps=30)
```

Guidance Scale	Time (it/s)
1.0	7.60
3.0	4.18
5.0	4.20
7.5	4.14
10.0	4.13
15.0	4.15



CLIP Score vs Guidance Scale for: "A professional photograph of a delicious pizza with pepperoni and cheese"

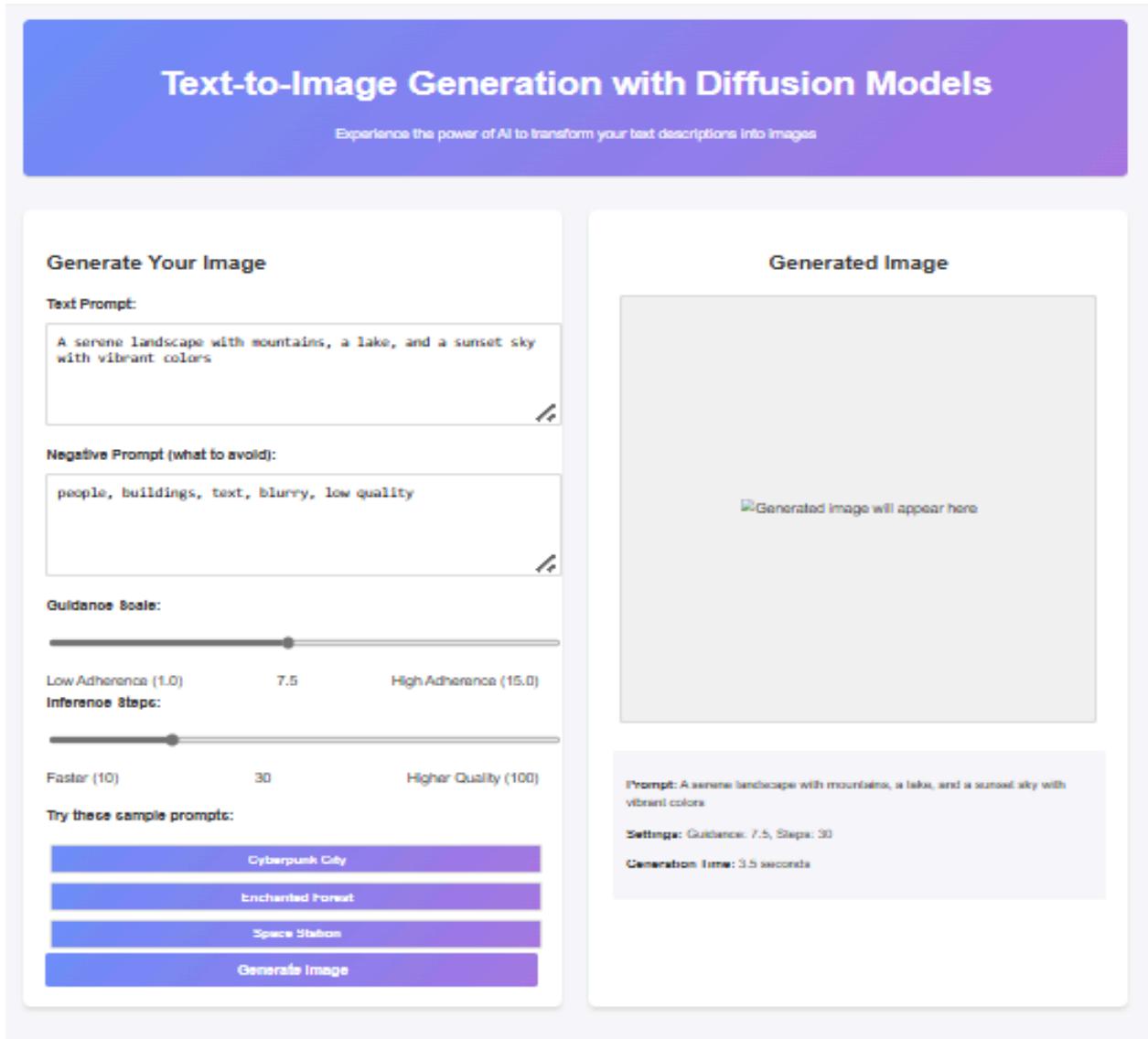


Key Findings and Observations

Through these experiments, several important findings were observed:

1. **Text Prompt Impact:** The text prompt significantly influences the generated image, with more detailed prompts generally producing more specific results.
2. **Guidance Scale Effect:** Higher guidance scale values lead to images that more closely follow the text prompt but may sacrifice some diversity and creativity.
3. **Inference Steps Trade-off:** More inference steps generally produce higher-quality images but require longer generation times.
4. **Negative Prompting Effectiveness:** Negative prompts effectively guide the model to avoid specific undesired elements in the generated images.
5. **CLIP Score Correlation:** CLIP scores generally increase with higher guidance scales, confirming better text-image alignment.

Web Interface Implementation



The project includes an interactive web interface that demonstrates the concept of text-to-image generation. The interface allows users to:

1. Enter text prompts describing desired images
2. Specify negative prompts for elements to avoid
3. Adjust guidance scale and inference steps
4. Choose from sample prompts
5. View the generated image and generation details

The web interface is implemented using HTML, CSS, and JavaScript, with the following key components:

```

<div class="app-container">
  <div class="input-section">
    <h2>Generate Your Image</h2>

    <label for="prompt">Text Prompt:</label>
    <textarea id="prompt" placeholder="Describe the image you want to generate...">A serene landscape with mountains, a lake, and a sunset sky with vibrant colors</textarea>

    <label for="negative-prompt">Negative Prompt (what to avoid):</label>
    <textarea id="negative-prompt" placeholder="Elements you want to avoid in the image...">people, buildings, text, blurry, low quality</textarea>

    <label for="guidance-scale">Guidance Scale:</label>
    <input type="range" id="guidance-scale" min="1" max="15" step="0.5" value="7.5">
    <div style="display: flex; justify-content: space-between;">
      <span>Low Adherence (1.0)</span>
      <span id="guidance-value">7.5</span>
      <span>High Adherence (15.0)</span>
    </div>

    <label for="steps">Inference Steps:</label>
    <input type="range" id="steps" min="10" max="100" step="5" value="30">
    <div style="display: flex; justify-content: space-between;">
      <span>Faster (10)</span>
      <span id="steps-value">30</span>
      <span>Higher Quality (100)</span>
    </div>

    <div class="sample-prompts">
      <p><strong>Try these sample prompts:</strong></p>
      <button class="sample-prompt-btn" data-prompt="A cyberpunk cityscape at night with neon lights and flying cars">Cyberpunk City</button>
      <button class="sample-prompt-btn" data-prompt="An enchanted forest with magical creatures and glowing plants">Enchanted Forest</button>
      <button class="sample-prompt-btn" data-prompt="A futuristic space station orbiting a ringed planet">Space Station</button>
    </div>
  </div>

```

Interface Features

- Text Prompt Input:** A text area for entering detailed image descriptions
- Negative Prompt Input:** A text area for specifying elements to avoid
- Guidance Scale Slider:** Controls how closely the generation follows the text
- Inference Steps Slider:** Controls the number of denoising steps
- Sample Prompts:** Pre-defined prompts for easy experimentation
- Image Display:** Area to view the generated image
- Generation Information:** Details about the generation process

Educational Components

The interface includes educational elements that explain:

- How diffusion models work
- The step-by-step denoising process
- The function of different parameters
- Tips for better results

Future Improvements

The project could be enhanced with the following improvements:

- API Integration:** Connect the web interface to a real Stable Diffusion API
- Additional Control Methods:** Implement inpainting, outpainting, and img2img capabilities
- Custom Fine-tuning:** Add the ability to fine-tune models on custom datasets

4. **Style Transfer:** Implement style transfer capabilities using text prompts
5. **Resolution Control:** Add options for different output resolutions

Ethical Considerations

The project acknowledges several important ethical considerations:

1. **Potential Misuse:** The technology could be misused for deepfakes or misinformation
2. **Copyright and Intellectual Property:** Questions about ownership of AI-generated content
3. **Bias and Representation:** Models may perpetuate or amplify biases in training data
4. **Content Moderation:** Importance of preventing harmful content generation

Conclusion

This project demonstrates the power of diffusion models for text-to-image generation, providing both theoretical understanding and practical implementation. Through the Jupyter notebook and web interface, users can explore how diffusion models transform text descriptions into high-quality images and experiment with different parameters to control the generation process.

The combination of mathematical foundations, neural architecture insights, and interactive demonstrations provides a comprehensive understanding of this cutting-edge technology and its applications in creative and practical domains.

References

1. Ho, J., Jain, A., & Abbeel, P. (2020). Denoising Diffusion Probabilistic Models. *Advances in Neural Information Processing Systems*.
2. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-Resolution Image Synthesis with Latent Diffusion Models. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
3. Hugging Face Diffusers Documentation. <https://huggingface.co/docs/diffusers/>
4. Stable Diffusion - CompVis, Stability AI, and LAION. <https://stability.ai/stable-diffusion>