
Fine-Tuning Stable Diffusion XL for Naruto-Style Image Generation with DreamBooth and LoRA

1. Project Overview

This project aims to generate high-quality Naruto-themed anime images from natural language prompts (e.g., "Naruto fighting in the rain") using Stable Diffusion XL (SDXL). By fine-tuning SDXL with DreamBooth and Low-Rank Adaptation (LoRA), we adapt the model to produce Naruto-style visuals. The process leverages the Hugging Face diffusers library and a custom Naruto dataset.

2. Background and Key Concepts

2.1 Diffusion Models

Diffusion models generate images by iteratively denoising random noise into meaningful visuals. Stable Diffusion, a latent diffusion model, operates in a compressed latent space to reduce computational costs.

- **How It Works:**
 - A forward process adds Gaussian noise to an image.
 - A reverse process, powered by a U-Net, predicts and removes noise to reconstruct the image.
 - A Variational Autoencoder (VAE) encodes images into latent space and decodes them back to pixels.
 - Text conditioning is achieved via a frozen CLIP text encoder, enabling classifier-free guidance.
- **Components:**
 - VAE (encoder/decoder)
 - U-Net (diffusion model)
 - CLIP text encoder
- **Advantages:**
 - High-fidelity image generation
 - Fine-grained text conditioning

2.2 Stable Diffusion XL (SDXL)

SDXL is an advanced version of Stable Diffusion with enhanced capabilities.

- **Key Features:**
 - 3.5 billion parameters
 - Larger U-Net (~3× compared to SD2.x)
 - Dual text encoders: OpenCLIP ViT-bigG/14 and original CLIP
 - Size-and-crop conditioning
 - Two-stage pipeline (base + refiner model)
 - Outputs high-resolution images (e.g., 1024×1024) with photorealistic details

2.3 DreamBooth

DreamBooth personalizes diffusion models for specific subjects or styles using 3–5 images.

- **Core Idea:**
 - Introduces a new token (e.g., [V]) associated with the target subject/style.
 - Enables prompts like “[V] in a forest” to generate context-specific images.
- **Training:**
 - Uses a prompt template (e.g., “a photo of a [V] dog”).
 - Employs prior-preservation loss with class images to prevent overfitting.
 - Phases: Low-resolution token binding, high-resolution refinement.

2.4 Low-Rank Adaptation (LoRA)

LoRA is a lightweight fine-tuning method that modifies a small subset of model parameters.

- **How It Works:**
 - Freezes original weights (W).
 - Adds trainable low-rank matrices (A , B) such that $W' = W + B \times A$.
 - Targets cross-attention layers for text-image alignment.
 - **Benefits:**
 - Reduces VRAM usage
 - Speeds up training
 - Maintains full fine-tuning quality
 - Simplifies deployment
-

3. Dataset

3.1 Description

We used the [lambdalabs/naruto-blip-captions](#) dataset from Hugging Face:

- **Content:** 1,221 Naruto-style anime images (~501 MB) with BLIP-generated captions.
- **Format:** Each entry includes a JPEG image and a text caption (e.g., “Naruto Uzumaki with orange background”).
- **Source:** Images from [Narutopedia](#), captioned using BLIP.

3.2 Preprocessing

- **Images:**
 - Resized to 256x256 and center-cropped if needed.
 - Converted to RGB and normalized to the SDXL latent space.
 - No augmentations applied to preserve identity learning.
 - **Captions:**
 - Used directly as prompts without cleaning.
 - Tokenized with SDXL’s CLIP tokenizer.
 - Appended to a fixed instance prompt: “a naruto anime character.”
 - **Special Tokens:**
 - No new tokens introduced; the class prompt “a naruto anime character” was used.
 - **Splitting:**
 - All images used for training (no validation split) due to the small dataset size.
 - DreamBooth’s regularization mitigates overfitting.
-

4. Fine-Tuning Setup

4.1 Model Architecture

Component	Details
Base Model	stabilityai/stable-diffusion-xl-base-1.0
VAE	madebyollin/sdxl-vae-fp16-fix
Text Encoder	Dual CLIP encoders (OpenCLIP + original)
Scheduler	DDIM with multi-step denoising
Fine-Tuning Method	DreamBooth + LoRA (attention layers only)

4.2 Training Configuration

Parameter	Value
Pretrained Model	stabilityai/stable-diffusion-xl-base-1.0
VAE Model	madebyollin/sd-xl-vae-fp16-fix
Dataset	lambdalabs/naruto-blip-captions
Instance Prompt	“a naruto anime character”
Resolution	256x256
Train Batch Size	1
Gradient Accumulation	4 (effective batch size = 4)
Learning Rate	1.0 (with Prodigy optimizer)
Text Encoder LR	1.0
Optimizer	Prodigy
SNR Gamma	5.0
LR Scheduler	Constant
Max Train Steps	500
Gradient Checkpointing	Enabled
Mixed Precision	fp16
LoRA Rank	4

Push to Hub

Enabled

4.3 Training Process

- **Hardware:** Single Tesla T4 GPU.
- **Duration:** ~1 hour with mixed precision.
- **Steps:**
 - Installed dependencies (diffusers, transformers, peft, accelerate, prodigyopt).
 - Downloaded Hugging Face's `train_dreambooth_lora_sd1x_advanced.py` script.
 - Configured accelerate for GPU and fp16.
 - Authenticated with Hugging Face Hub using a token.
 - Launched training with accelerate launch, specifying hyperparameters.
 - Trained for 500 diffusion steps, updating LoRA parameters in U-Net's cross-attention layers.
 - Pushed LoRA weights to Hugging Face Hub.
- **Training Dynamics:**
 - U-Net processes batches of latents + noise and text embeddings.
 - Computes MSE loss on predicted noise.
 - Prodigy optimizer updates LoRA matrices.
 - SNR weighting ($\gamma=5.0$) emphasizes low-noise timesteps.
 - Gradient checkpointing reduces VRAM usage.

4.4 Attention Mechanisms

- **SDXL's Cross-Attention:**
 - Aligns image latents (queries) with text embeddings (keys/values) from dual CLIP encoders.
 - Enables spatial feature binding to prompt words (e.g., "Naruto" → orange hair).
 - **LoRA's Role:**
 - Adds low-rank matrices (A, B) to query/value projections in cross-attention layers.
 - Rank-4 updates perturb weights minimally, capturing Naruto-style features.
 - Post-training, LoRA matrices are fused into U-Net weights for efficient inference.
-

5. Code Walkthrough

The fine-tuning was implemented in a Jupyter notebook:

1. **Setup:**
 - Installed libraries and downloaded the training script.
 - Configured accelerate and authenticated with Hugging Face.

2. Training:

- Executed !accelerate launch with parameters (see Section 4.2).
- Script loaded SDXL, applied LoRA, and processed the dataset with the instance prompt.

3. Post-Training:

- Cleared GPU memory.
- Loaded the fine-tuned pipeline with StableDiffusionXLPipeline and fused LoRA weights.

Inference Example:

python

Copy

```
prompt = "a naruto anime character with red hair and green eyes"
negative_prompt = "low quality, worst quality, bad anatomy, bad composition"
image = pipe(
    prompt=prompt,
    negative_prompt=negative_prompt,
    guidance_scale=7.5,
    num_inference_steps=30
```

4.).images[0]

- Generated a 1024×1024 Naruto-style image for visual inspection.

6. Evaluation

6.1 Qualitative Results

● Prompts Tested:

- "Naruto fighting in the forest"
- "A Naruto anime-style portrait"
- "Sasuke Uchiha with red Sharingan"

● Observations:

- Images exhibit Naruto's anime aesthetic (bright colors, sharp outlines, ninja headbands).
- Negative prompts reduce artifacts like bad anatomy.

6.2 Quantitative Metrics

- **Fréchet Inception Distance (FID):** Measures similarity between generated and real Naruto images (lower is better).

- **CLIP Score:** Evaluates text-image alignment via cosine similarity of CLIP embeddings (higher is better).
 - **Limitations:** Metrics may be less reliable for narrow styles; human judgment is critical.
-

7. MLOps Considerations

To productionize the pipeline:

Tool	Purpose
MLflow	Track experiments, parameters, checkpoints
TensorBoard	Monitor loss and learning rate
Hugging Face Hub	Version and deploy models
Docker	Ensure reproducible environments
FastAPI	Create a prompt-to-image web API

- **CI/CD:** Use GitHub Actions for automated testing and small-scale training.
- **Serving:** Deploy via Docker/Kubernetes with NVIDIA Triton or Hugging Face Inference API.
- **Monitoring:** Track latency, throughput, and output quality.
- **Data Management:** Version datasets with DVC or Delta tables.

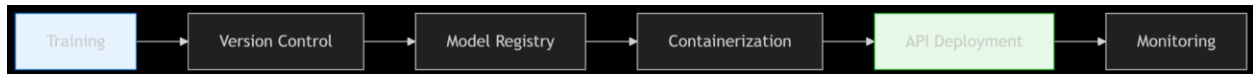
8. Challenges and Limitations

- **Overfitting:** Small dataset (~1.2K images) risks memorization. Mitigated by DreamBooth's prior-preservation and limited steps.
 - **Caption Noise:** BLIP captions may be inconsistent (e.g., unrelated scenes).
 - **Generalization:** Prompts outside the Naruto domain (e.g., "Naruto in space") may yield poor results.
 - **Memory:** LoRA reduces VRAM needs, but SDXL is still resource-intensive.
-

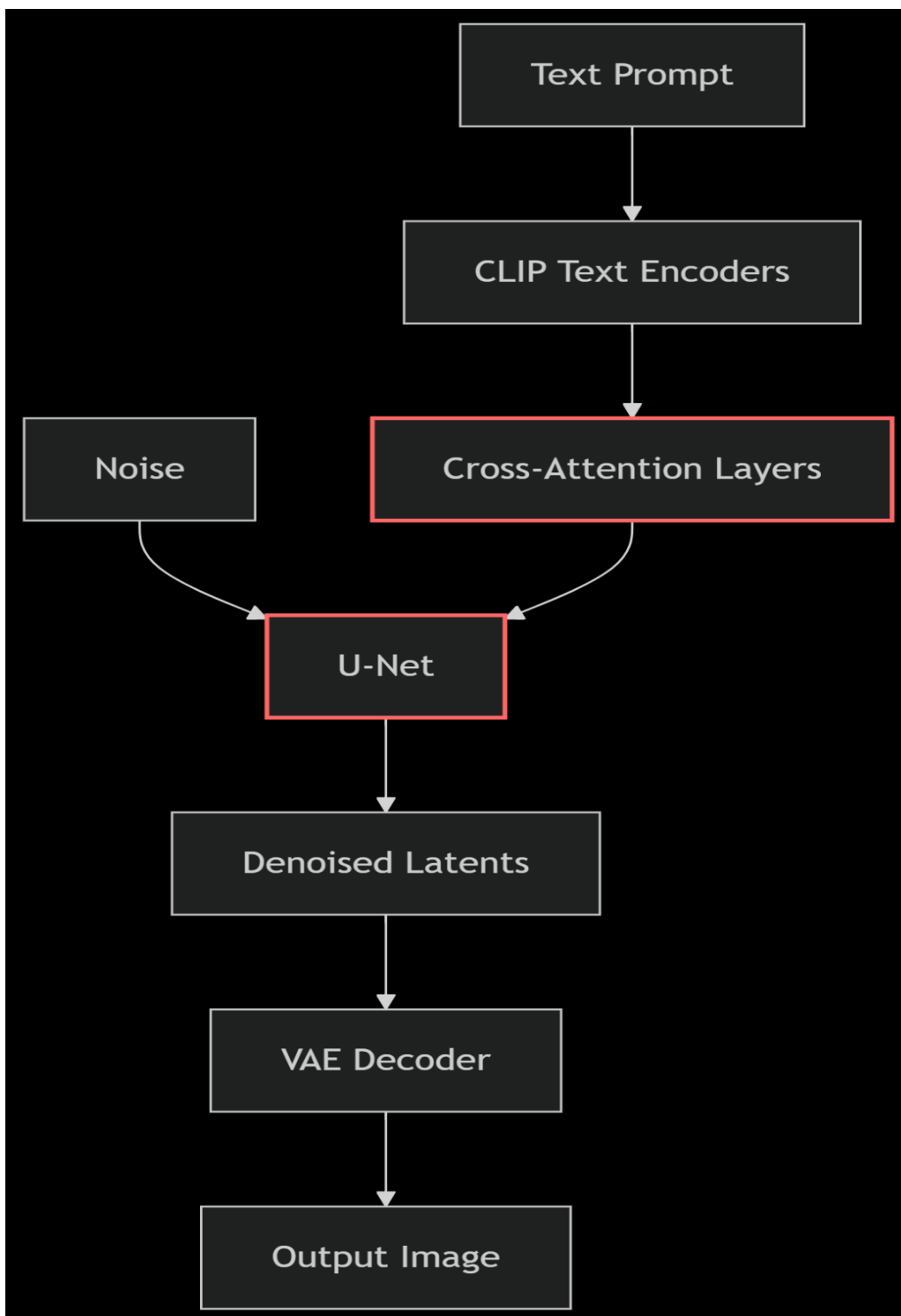
9. Future Work

- **Dataset Expansion:**
 - Curate cleaner captions or re-caption with specialized models.
 - Include multi-character images for complex prompts (e.g., “Naruto and Sasuke”).
 - **Prompt Engineering:**
 - Introduce unique tokens via textual inversion (e.g., <naruto-style>).
 - Combine LoRA with textual inversion for finer control.
 - **Multi-Style Training:**
 - Extend to other anime styles (e.g., One Piece) using multi-token or multi-head LoRA.
 - **Advanced Techniques:**
 - Apply dropout, weight decay, or augmentations to reduce overfitting.
 - Explore gradual unfreezing or style transfer objectives.
-

10.MLOps Pipeline



11.Model Architecture



12. Conclusion

This project successfully fine-tuned Stable Diffusion XL to generate Naruto-style anime images using DreamBooth and LoRA. Key achievements:

- Efficient fine-tuning with minimal parameters via LoRA.
- High-quality Naruto-themed outputs aligned with input prompts.
- Practical MLOps strategies for production deployment.

The approach is extensible to other styles and demonstrates the power of combining advanced diffusion models with lightweight fine-tuning techniques.