

LCM-LoRA Model Report

Model Overview	3
How does it work?	3
Consistency Models.....	3
Latent Consistency Model (LCM)	4
LCM-LoRA.....	5
Sampling Method for LCM-LoRA.....	5
Code Overview.....	6
Cell 1 Downloading dependencies	6
Cell 2 Building and loading model	6
Cell 3 Generating Images	7
Results	8

Model Overview

To understand this model, we will have to delve into the realm of Consistency Models (CM), a novel category of diffusion models engineered for single-step image generation.

Originating from the work of Yang Song and colleagues in their publication on "Consistency Models" these models represent a breakthrough in efficient image synthesis.

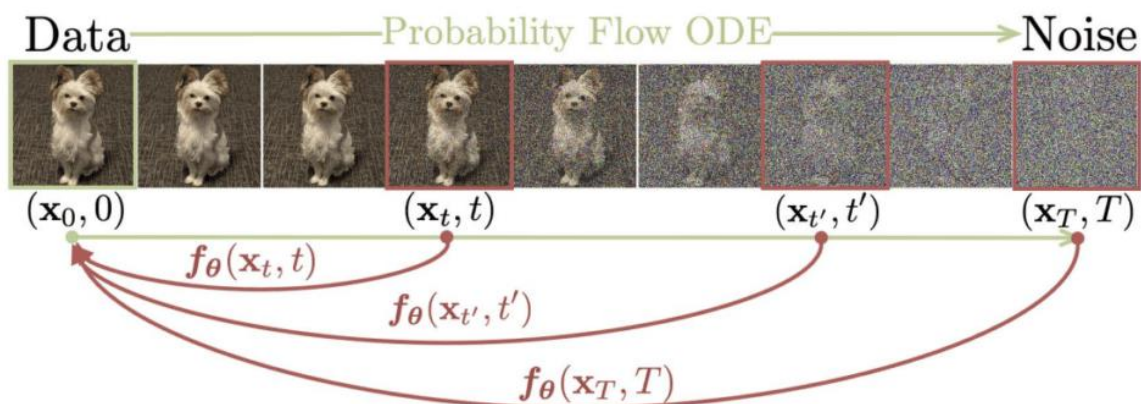
The Latent Consistency Model (LCM) extends this concept to latent diffusion models, notably Stable Diffusion, wherein image denoising occurs within the latent space. Typically, employing LCM necessitates training a new model for each custom checkpoint, which can be cumbersome.

How does it work?

To comprehend the workings of LCM-LoRA, it's imperative to grasp the concept of consistency models and the endeavors aimed at enhancing diffusion models' speed.

Consistency Models

A consistency model is a type of diffusion model trained to generate an AI image in a single step. Acting as a more efficient student model, it undergoes training with a teacher model, such as the SDXL model. The objective is for the student model to produce the same image as the teacher model, albeit in a single step. Essentially, the consistency model serves as a swifter alternative to the teacher model.



The concept behind the consistency model revolves around establishing a correlation between the final AI image and any denoising step within a diffusion model's progression.

For instance, in a scenario where a diffusion model is trained to generate an AI image across 50 steps, a consistency model maps the intermediate noisy images at each step (0, 1, 2, 3...) to the final step (50).

Named for its emphasis on consistency, the training of a consistency model capitalizes on the reliability of the mapping process: Each intermediate image is consistently mapped to the final image within a single step. This means that regardless of the noise level in the image, the mapping function consistently produces outputs identical to the final image.

In practical application, the image quality from a single-step generation is typically subpar. Therefore, practitioners often opt for a few steps instead.

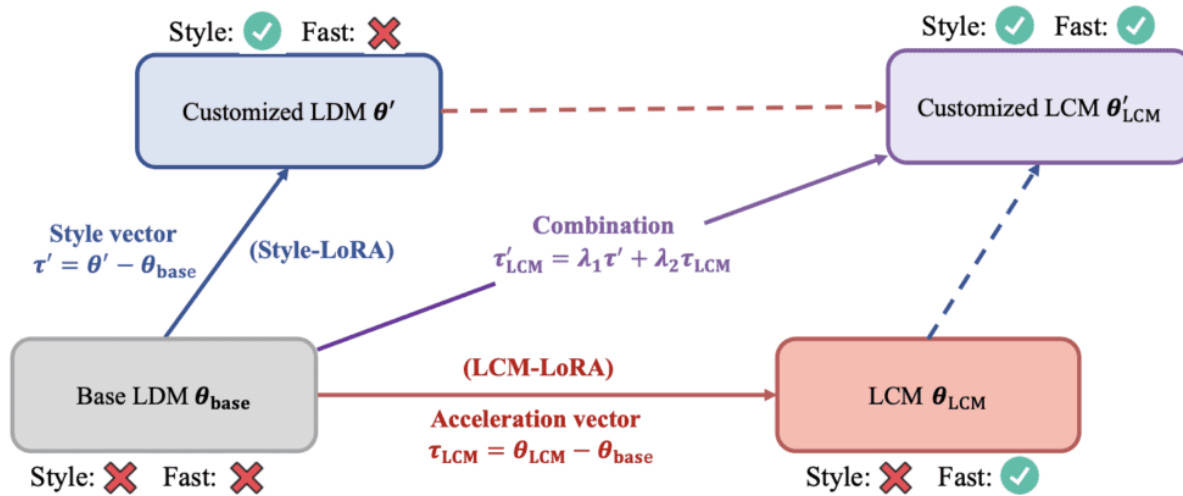
Functioning as a distillation method, the consistency model leverages information extraction and reorganization from an existing (teacher) model to enhance efficiency.

Compared to the progressive distillation method, which pledged significant speed improvements for Stable Diffusion, consistency models demonstrate superiority by generating higher-quality images.

Latent Consistency Model (LCM)

The Latent Consistency Model (LCM) represents a form of consistency model integrated with latent diffusion techniques, as observed in Stable Diffusion. Explored by Simian Luo and colleagues in their paper "Latent Consistency Models: Synthesizing High-Resolution Images with Few-Step Inference," the distinguishing factor of LCM lies in its operation within the latent space, contrasting with the original consistency model, which operates in pixel space.

LCM-LoRA



Instead of training a new checkpoint model, you utilize LoRA to modify an existing LCM. LoRA functions like a compact patch, offering several advantages:

- Portability:

LCM-LoRA can be applied to ANY Stable Diffusion checkpoint models, including v1.5 and SDXL models, enhancing their speed.

- Faster training:

With fewer weights to train, LoRA accelerates the training process and reduces resource demands.

Sampling Method for LCM-LoRA

The LCM model aims for 1-step inference to generate the final AI image swiftly. However, this approach may not always deliver optimal quality. Here's how the sampling method for LCM works:

1. Denoise the latent image.
2. Introduce noise back into the image based on the noise schedule.
3. Repeat steps 1 and 2 until reaching the final sampling step.

Code Overview

Cell 1 Downloading dependencies

```
!pip install --quiet --upgrade diffusers accelerate
!pip install --quiet -U peft transformers
```

The code consists of two commands written in the Markdown syntax for installing Python packages using the pip package manager. The first command installs the *diffusers* and *accelerate* packages, and the second command installs the *peft* and *transformers* packages. The *--quiet* flag is used in both commands to suppress the output of the installation process, and the *--upgrade* flag in the first command and the *-U* flag in the second command ensure that the latest version of each package is installed. These packages are likely needed for machine learning or data analysis tasks in a Python environment.

Cell 2 Building and loading model

```
SD_Models = "runwayml/stable-diffusion-v1-5"
LoRA_Models = "latent-consistency/lcm-lora-sdv1-5"

from diffusers import DiffusionPipeline, LCMScheduler
import torch
import time

pipe = DiffusionPipeline.from_pretrained(
    SD_Models,
    torch_dtype=torch.float16,
    use_safetensors=True,
    variant="fp16",
)

pipe.load_lora_weights(LoRA_Models)
pipe.scheduler = LCMScheduler.from_config(pipe.scheduler.config)

pipe = pipe.to("cuda", dtype=torch.float16)
```

The code initializes a diffusion pipeline using the *DiffusionPipeline* class from the *diffusers* library and loads the pre-trained weights for a stable diffusion model and a partially custom-trained LoRA model. The pipeline is configured to use half-precision floating-point and the *LCMScheduler* for more efficient computations.

The LoRA model is a custom-trained model on the images of cracks you provided using the *kohya_ss* software, a third-party tool that I used to train model locally on my personal computer. The LoRA model is designed to fine-tune the pre-trained weights of a diffusion model, allowing me to customize the model's behavior for specific applications or domains.

After loading the pre-trained weights and configuring the pipeline, the code moves the pipeline to the GPU and sets the data type to half-precision floating-point for more efficient computation. The pipeline is then ready to be used for generating new images or other outputs based on the customized behavior of the LoRA model.

Overall, the code sets up a diffusion pipeline that incorporates the customized behavior of the LoRa model into the diffusion process. The pipeline is initialized and configured with the necessary components and parameters for efficient and effective computation.

Cell 3 Generating Images

```
Prompt = "A single spilt crack on a concrete wall " # @param {type:'
No_of_Steps = 4 # @param {type:"integer"}
Guidance_Scale = 1 # @param {type:"number"}

start_time = time.time()
images = pipe(
    prompt=Prompt,
    num_inference_steps=No_of_Steps,
    guidance_scale=Guidance_Scale,
).images[0]
end_time = time.time()

total_time = end_time - start_time
print(f"Total time taken for Generating Image: {total_time} seconds'
images
```

The code provided is a form that prompts you to enter a text prompt, the number of inference steps, and the guidance scale. Your inputs are stored in the *Prompt*, *No_of_Steps*, and *Guidance_Scale* variables, respectively.

The code then initializes a timer using the *time* module and generates an image using the *pipe* object, which is an instance of the *DiffusionPipeline* class created in a previous code. The *pipe* object is called with the *prompt*, *num_inference_steps*, and *guidance_scale* parameters.

The generated image is then displayed using the `.images[0]` attribute, which returns the first generated image in the output. The timer is stopped, and the total time taken for generating the image is calculated and printed.

Overall, the code provides a user interface for generating images using the diffusion pipeline with customized parameters. Your inputs are used to generate a unique image, and the total time taken for the generation is recorded and displayed. This allows you to experiment with different prompts and parameters to generate high-quality images efficiently.

Results

I've noticed that the model performs exceptionally well when provided with accurate and descriptive prompts, achieving a 100% accuracy in image generation. Here are examples of prompts along with their corresponding outputs: Prompt: Sunlight shining on an outdoor concrete wall with split cracks

Generated Image:



Prompt: Vertically split cracks on an old yellow painted concrete wall

Generated image:



Prompt: Water draining down next to a cracked concrete wall

Generated image:



These generations are not always the first output of the prompt but after running the prompt 2-3 times you will get your desired output