



AI VIET NAM
@aivietnam.edu.vn

Emoji Image Generation Project

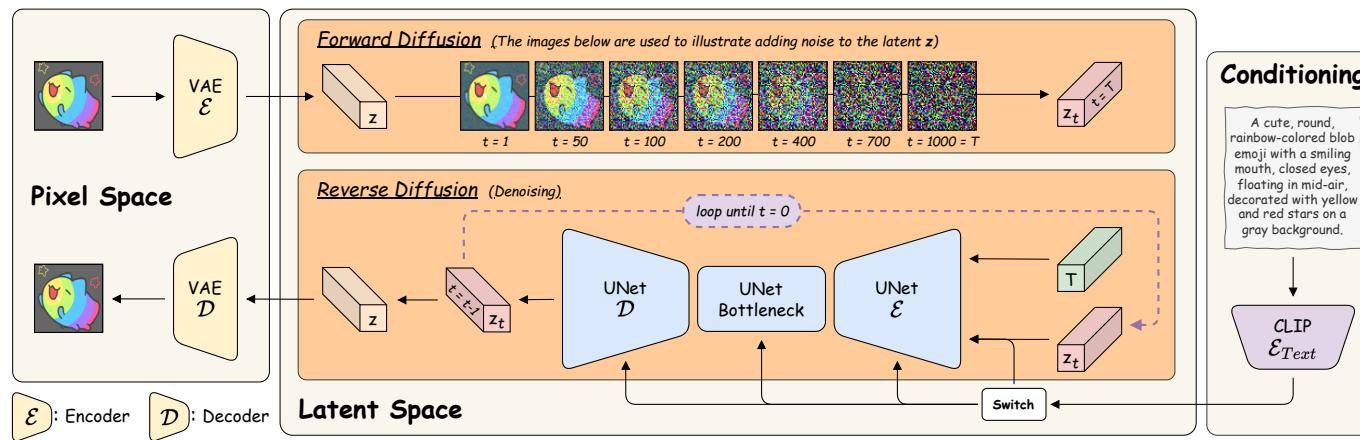
Dinh-Thang Duong – TA

Yen-Linh Vu – STA

Anh-Khoi Nguyen – STA

Getting Started

❖ Content



Our objectives:

- Investigate the task of Image Generation (IG) and Stable Diffusion (SD).
- Implement a Python program to crawl Discord Emojis on website.
- Implement and training a SD model from scratch using PyTorch and Diffusers.
- Implement and training a SD model from a pre-trained model using Diffusers.

Outline

- ❖ Introduction
- ❖ Data Collection and Preparation
- ❖ Training SD from scratch
- ❖ Training SD from pre-trained models
- ❖ Question



AI

AI VIET NAM
@aivietnam.edu.vn

Introduction

Introduction

❖ Getting Started

Forbes

The AI-Generated Studio Ghibli Trend, Explained

OpenAI's new image generator has released a surge of Studio Ghibli inspired images that many believe go against the philosophy of the legendary animation studio and director Hayao Miyazaki.

By [Dani Di Placido](#), Senior Contributor. Dani Di Placido covers film, television, and culture. Follow Author

Published Mar 27, 2025, 10:31am EDT, Updated Mar 31, 2025, 12:41pm EDT

Share Save Comment 6



Image created



Ask anything



Example: An image generation trend with Studio Ghibli generated image.

Introduction

❖ Introduction to the task of Image Generation

A yellow bird with brown and white wings and pointed bill



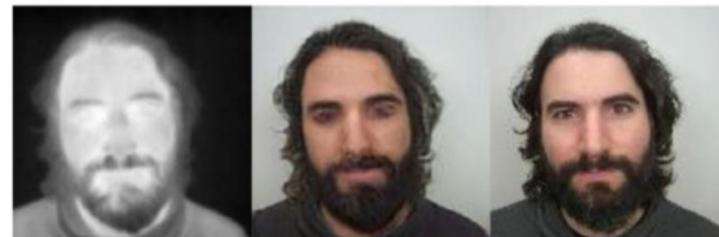
(a) Text-to-image generation MirrorGAN [55]



(c) Sketch-to-image generation EdgeGAN [97]



(e) Pose-Guided-image generation MsCGAN [94]



(g) Face generation [62]



(b) Scene-graph-to-image generation [26]



(d) Layout-to-image generation OC-GAN [74]



(f) Image-to-image Translation C2GAN [39]



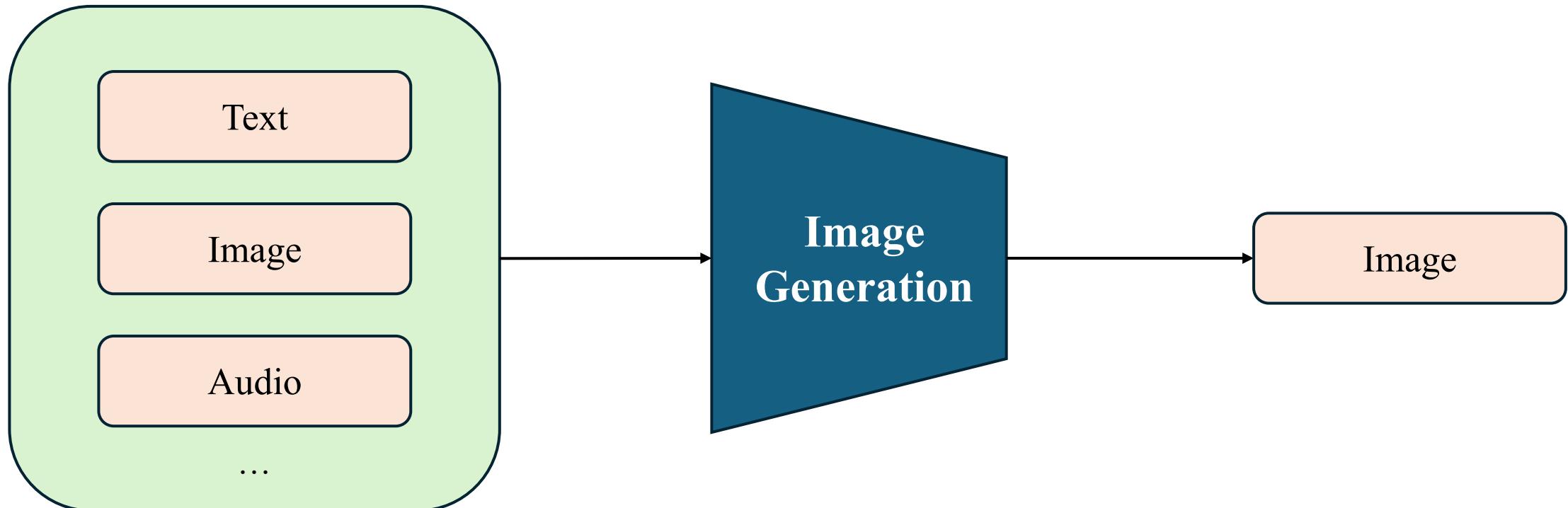
(h) Image generation PNAPGAN [89]

Image generation is the process of synthesizing new images from given inputs or learned representations.

It aims to produce visually coherent, contextually relevant, and high-quality outputs aligned with the intended goal.

Introduction

❖ Examples of Image Generation



Introduction

❖ Image Generation applications

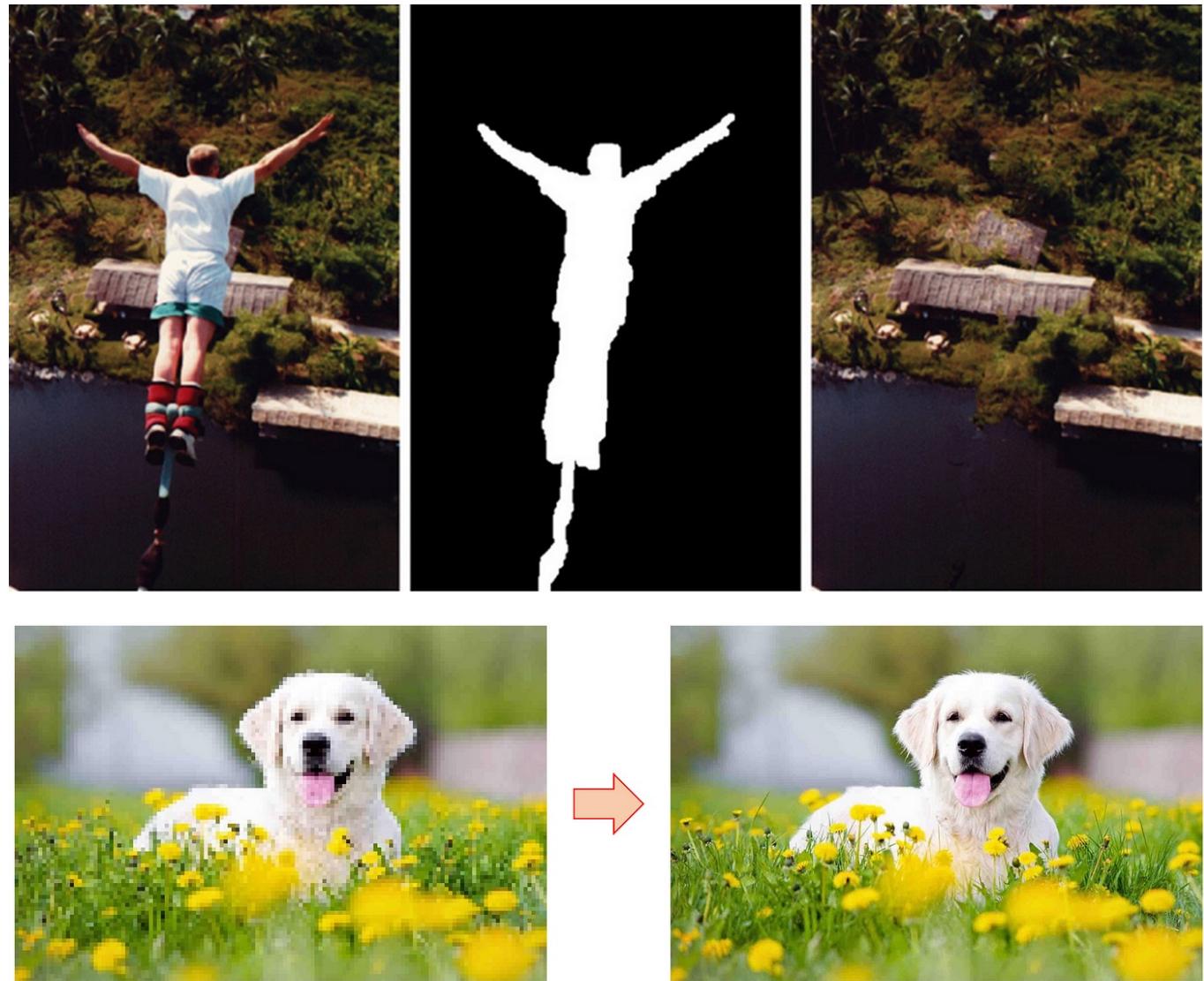
Inference Examples ⓘ

Text-to-Image Example 1

picture of a tiger, artstation

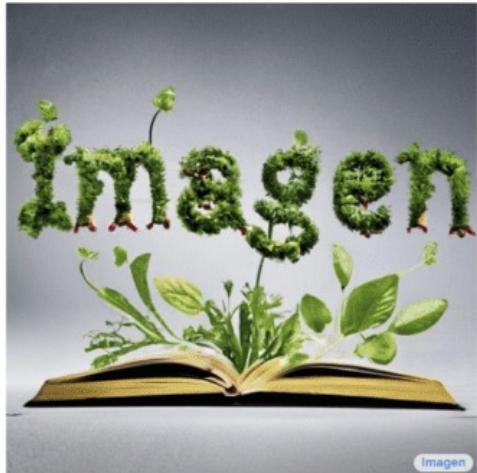


Maximize



Introduction

❖ Introduction to text2img



Sprouts in the shape of text 'Imagen' coming out of a fairytale book.



A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.



A high contrast portrait of a very happy fuzzy panda dressed as a chef in a high end kitchen making dough. There is a painting of flowers on the wall behind him.

Text-to-image generation refers to the task of automatically producing images that are semantically aligned with a given textual description or prompt.

By processing the linguistic content of the prompt and mapping it to visual features, models in this domain transform input text into synthesized images that reflect the conceptual meaning conveyed by the words.

Introduction

❖ Text2img challenges



- **Language Nuance:** Interpreting subtle, implied, or context-dependent details from textual descriptions.
- **Visual Accuracy:** Rendering accurate shapes, colors, and spatial relationships aligned with the prompt.
- **Realism and Style:** Achieving photorealism or adhering to a specific artistic style.



Introduction

❖ Text2img popular applications

Google DeepMind About Research Technologies Discover



Imagen

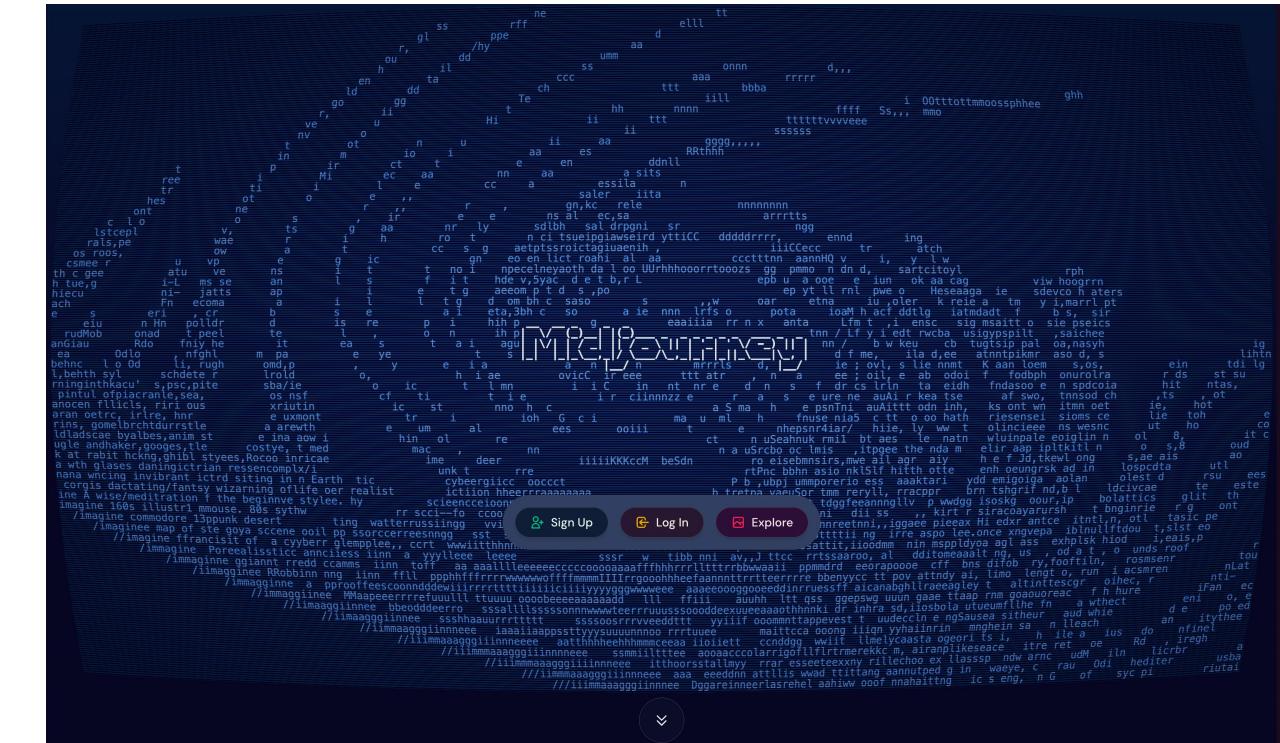
Our highest quality text-to-image model

Try in Gemini >

Try in ImageFX >



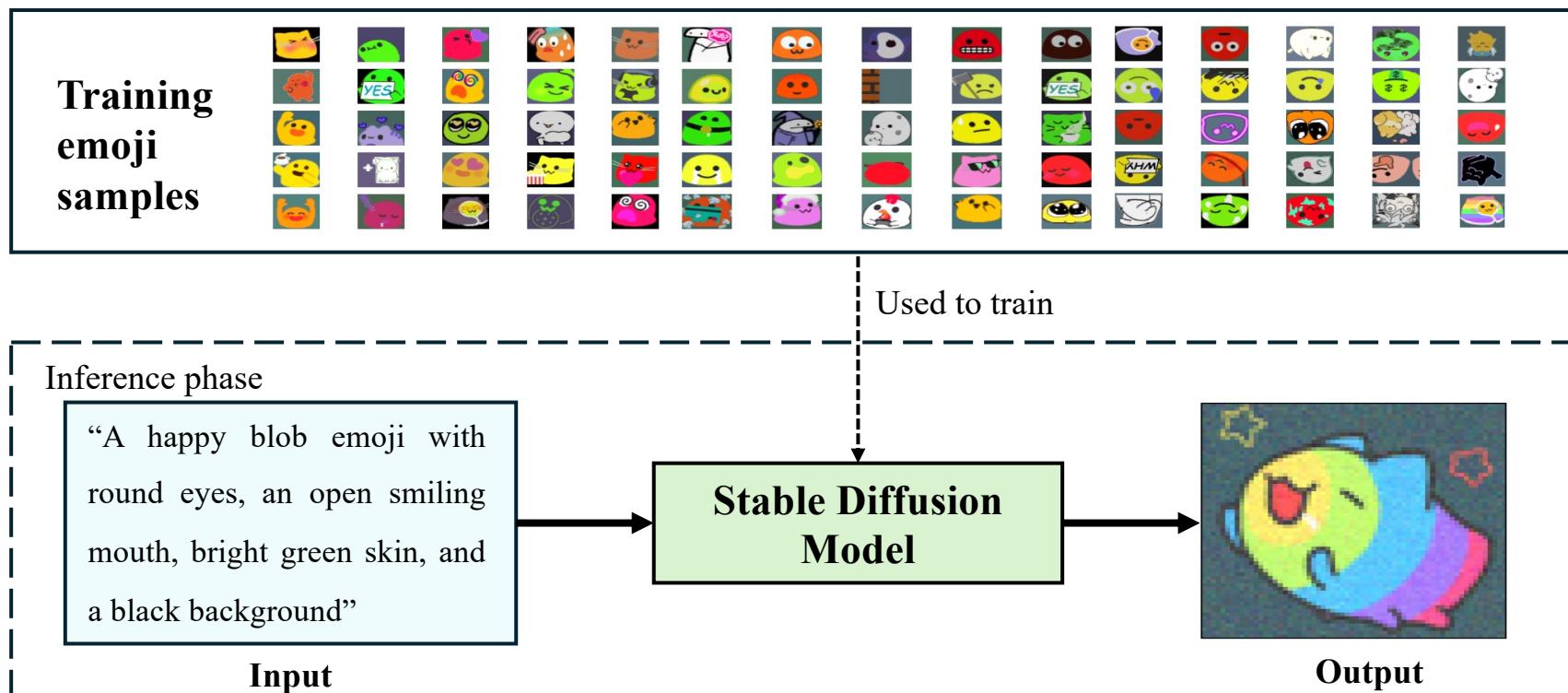
Search



Training SD from scratch

❖ Problem Statement

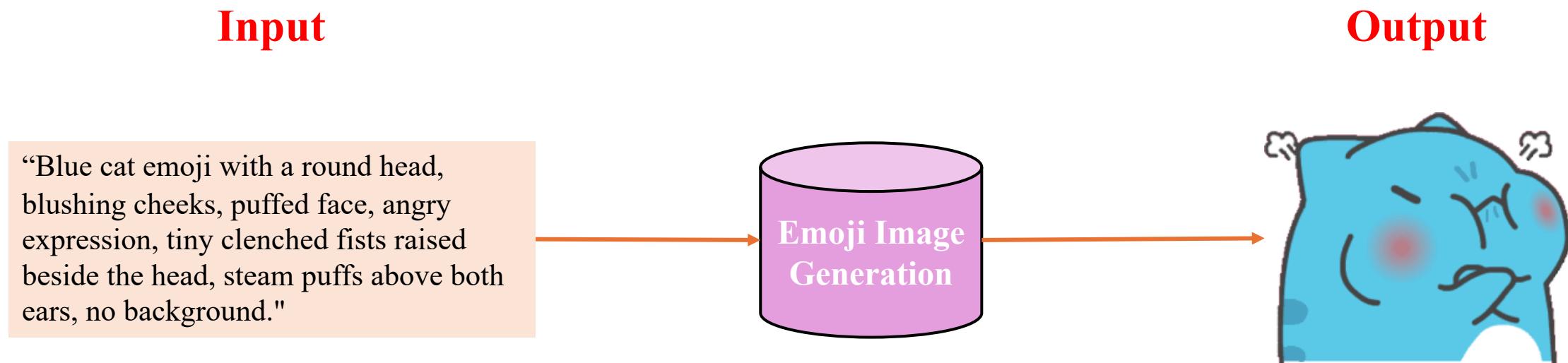
Problem Statement: Given a prompt that describes how the emoji looks like, generate a new emoji image using Stable Diffusion model.



Introduction

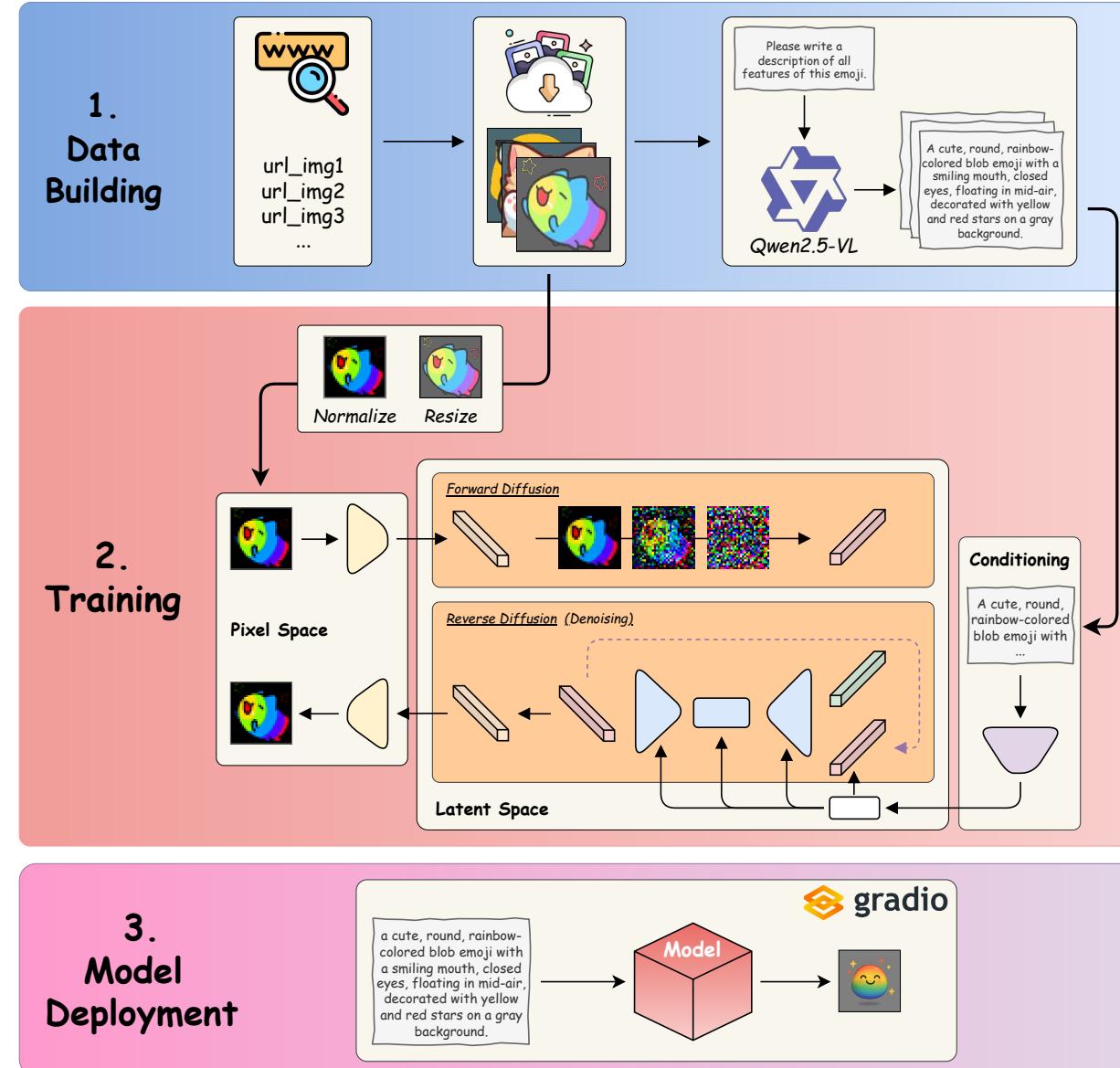
❖ Description

Problem Statement: Given a prompt that describes how the emoji looks like, generate a new emoji image using Stable Diffusion model.



Introduction

❖ Project Pipeline

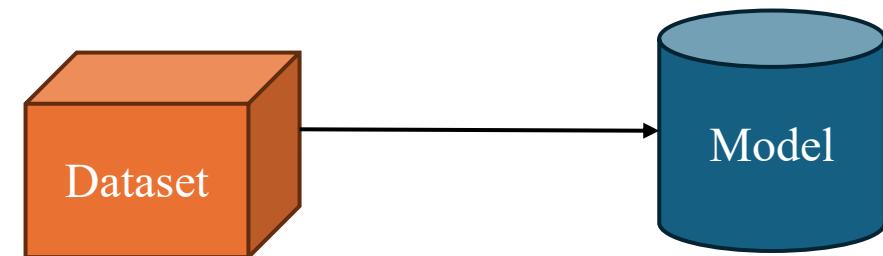


Data Collection and Preparation

Data Collection and Preparation

❖ Problem

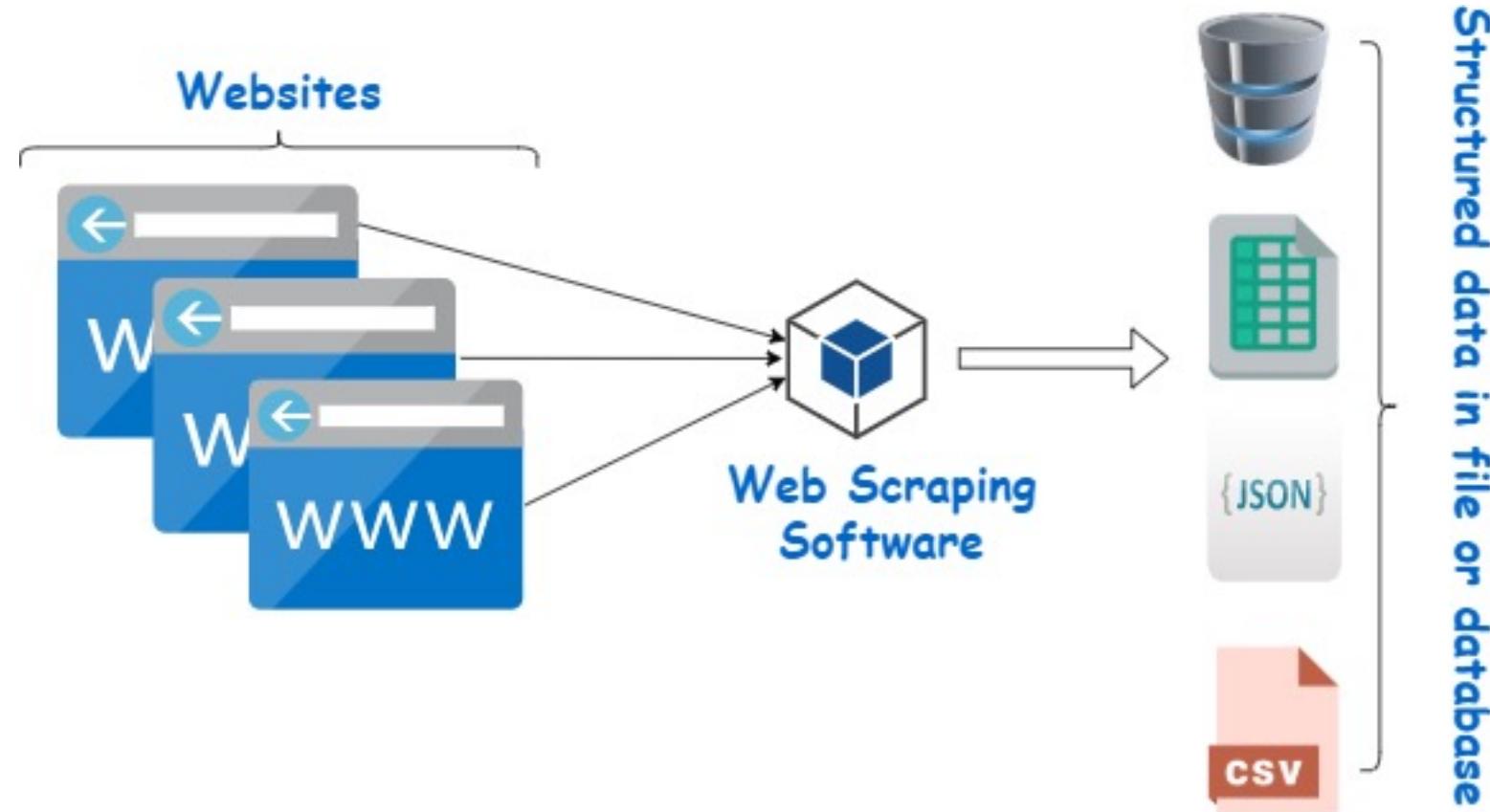
- ▼  blobs_crawled_data
- ▼  images
 -  Blobs_0000001.jpg
 -  Blobs_0000002.jpg
 -  Blobs_0000003.jpg
 -  Blobs_0000004.jpg
 -  Blobs_0000005.jpg
 -  Blobs_0000006.jpg
 -  Blobs_0000007.jpg
 -  Blobs_0000008.jpg
 -  Blobs_0000009.jpg
 -  Blobs_0000010.jpg



We don't have data yet.

Data Collection and Preparation

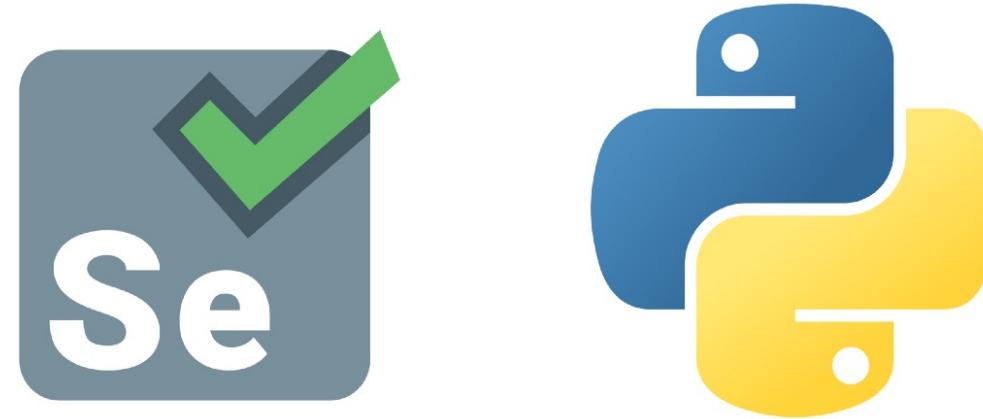
❖ Solution: Data Scraping



We can get data from websites.

Data Collection and Preparation

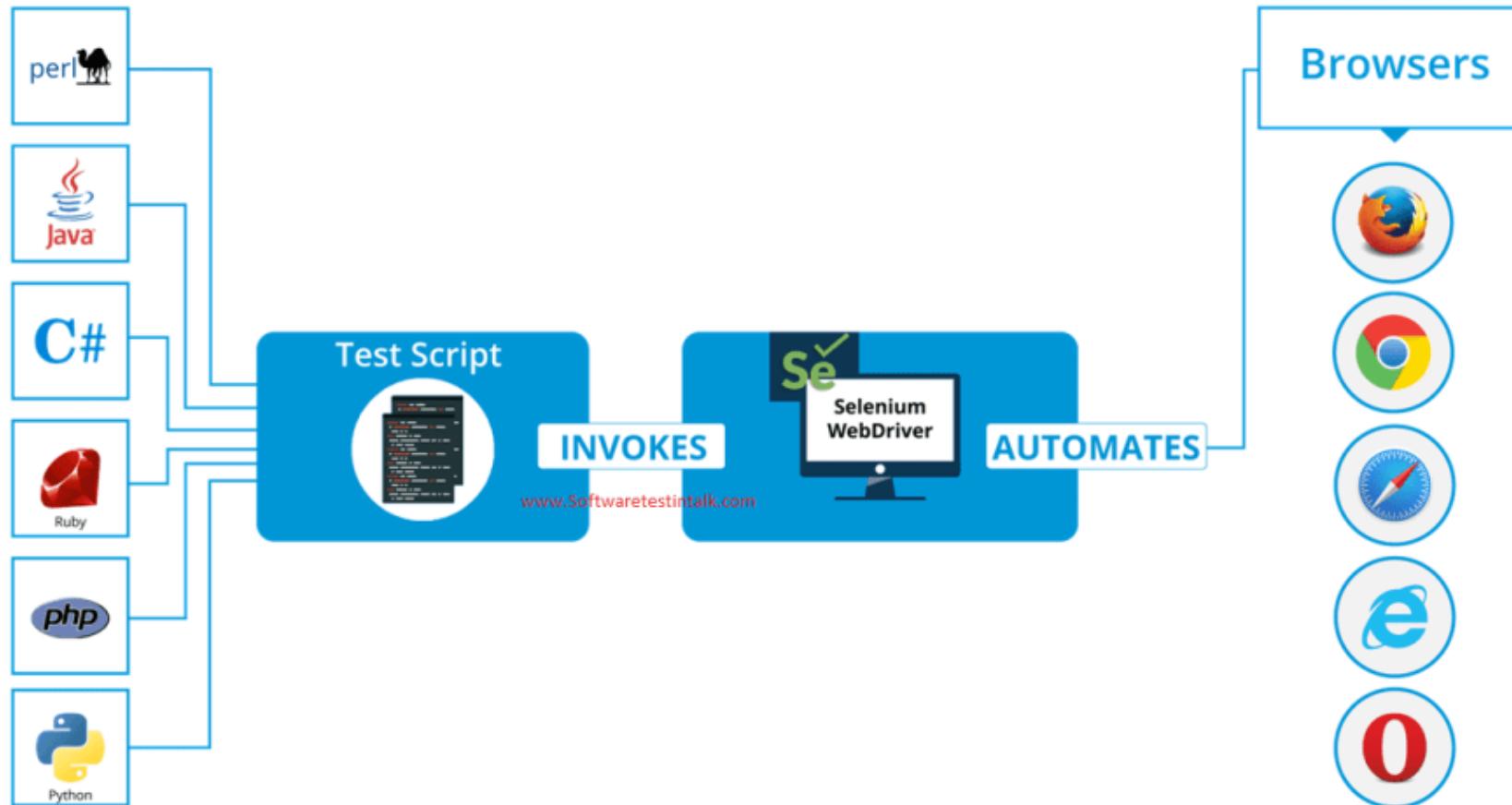
❖ Selenium Package



Selenium Package: Use to automate web browser interaction from Python.

Data Collection and Preparation

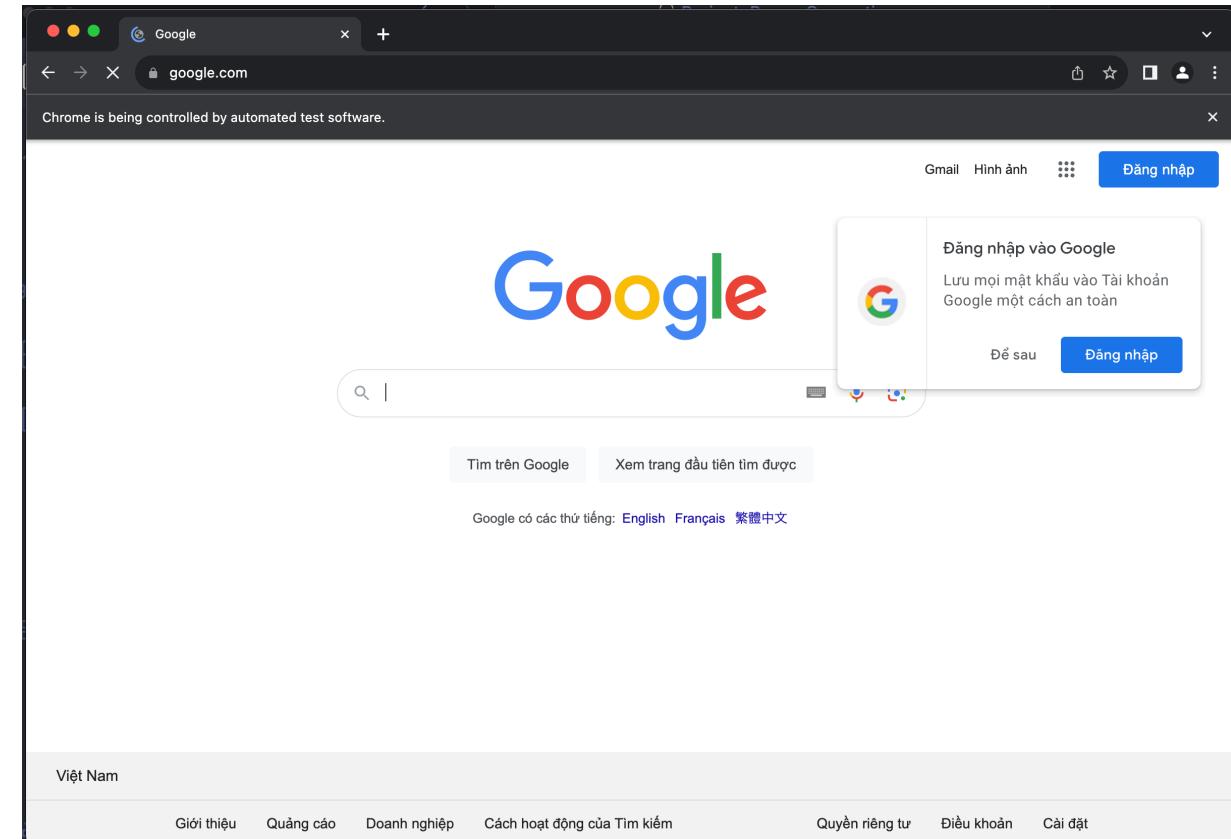
❖ Introduction to Selenium



Data Collection and Preparation

❖ Introduction to Selenium: Initialize driver

```
1 from selenium import webdriver
2
3 url = 'http://www.google.com'
4
5 driver = webdriver.Chrome()
6 driver.get(url)
```



Data Collection and Preparation

❖ Introduction to Selenium: HTML

```
<!DOCTYPE html>
<html>
    <!-- Root Element -->
    <head>
        <!-- Contains the header information -->
        <title>Title of the Page</title>
        <!-- Defines the title of the page -->
    </head>
    <body>
        <!-- Holds the content of the page -->
        <!-- Tags related to layout and formatting go here -->
    </body>
</html>
```

Data Collection and Preparation

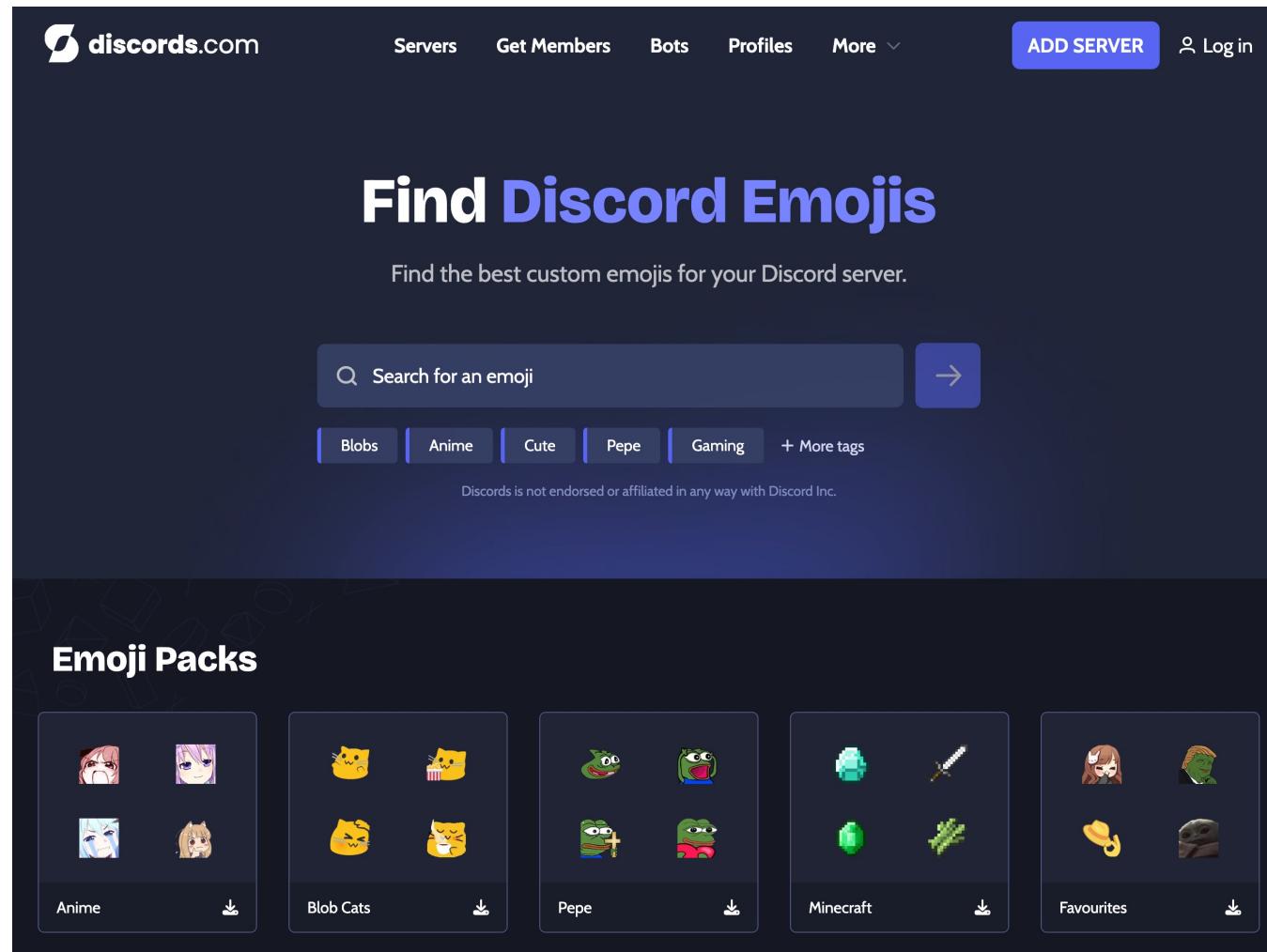
❖ Coding: Import libraries and initialize webdriver

```
1 import os
2 import requests
3 import time
4 import json
5
6 from tqdm import tqdm
7 from selenium import webdriver
8 from selenium.webdriver.chrome.service import Service
9 from selenium.webdriver.common.by import By
10 from selenium.webdriver.support.ui import WebDriverWait
11 from selenium.webdriver.support import expected_conditions as EC
```

```
1 WEBDRIVER_DELAY_TIME_INT = 10
2 TIMEOUT_INT = 10
3 service = Service(executable_path=r"/usr/bin/chromedriver")
4 chrome_options = webdriver.ChromeOptions()
5 chrome_options.add_argument("--headless")
6 chrome_options.add_argument("--no-sandbox")
7 chrome_options.add_argument("--disable-dev-shm-usage")
8 chrome_options.add_argument("window-size=1920x1080")
9 chrome_options.headless = True
10 driver = webdriver.Chrome(service=service, options=chrome_options)
11 driver.implicitly_wait(TIMEOUT_INT)
12 wait = WebDriverWait(driver, WEBDRIVER_DELAY_TIME_INT)
```

Data Collection and Preparation

❖ Website



Consider an emoji website: [Discord Emojis](#).

Data Collection and Preparation

❖ Website: Tags

Find Discord Emojis

Find the best custom emojis for your Discord server.

Search for an emoji →

Blobs Anime Cute Pepe Gaming – Less tags

Cat Panda Happy Sad Crying Laughing Crypto

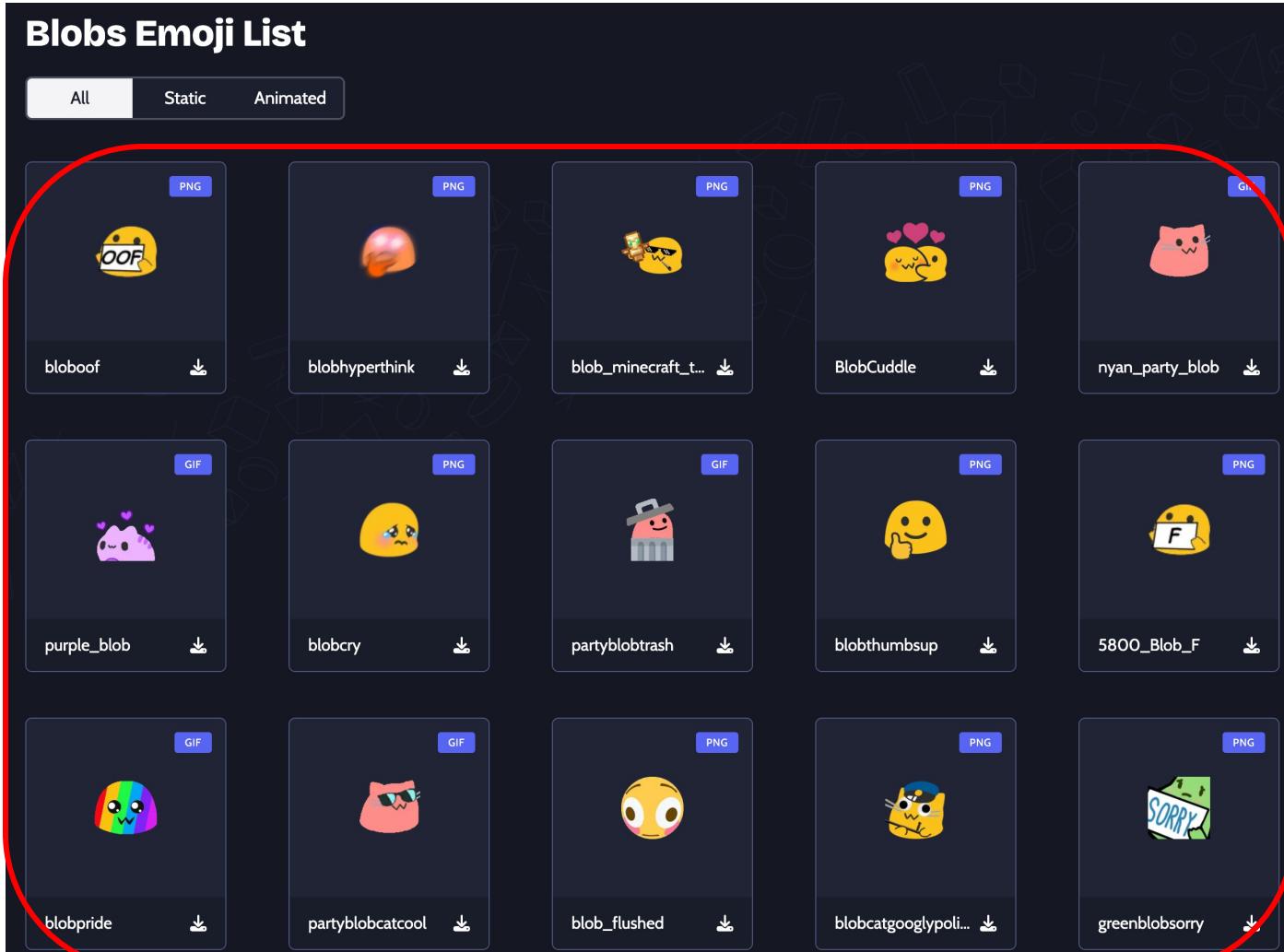
Furry Memes Thinking Purple Blue Red Animated



In our scope, **only Blobs** emojis are considered to be used for training the model.

Data Collection and Preparation

❖ Analyze website



A list of Blobs emojis are represented as div class “y-5H7YS5”. 25

Data Collection and Preparation

❖ Analyze website



The screenshot shows the browser's element inspector with the following details:

- HTML Structure:**

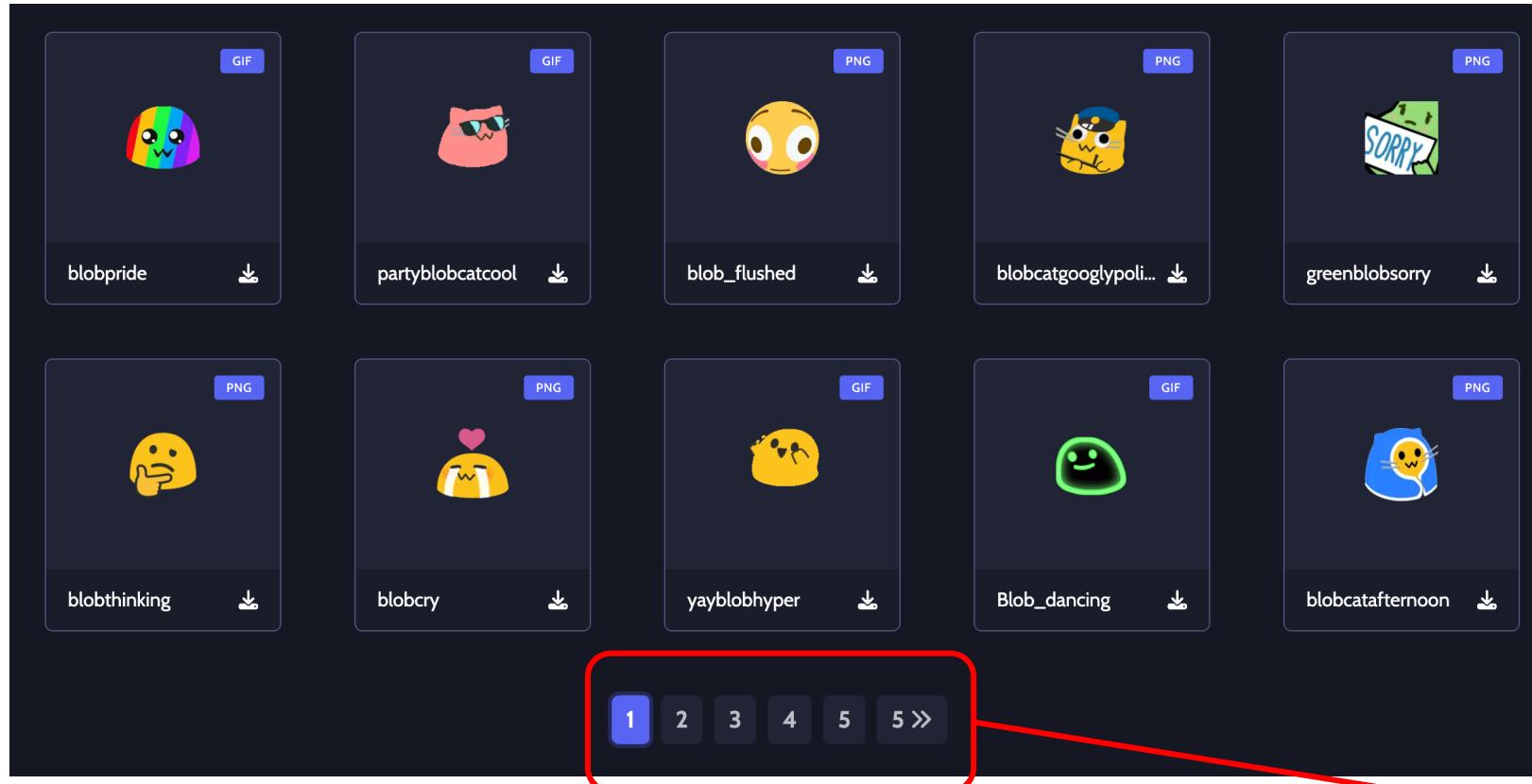
```
<div class="LQY5mtmC"> flex
  > <div class="G020KZik"> ... </div>
  > <div class="aLnnpRah text-center"> flex
    > <div class="Mw1EAtrx">
      
```
- Image Preview:** A yellow smiley face emoji with a thumbs up, displayed on a checkered background.
- Image Properties:**
 - Rendered size: 60 × 60 px
 - Rendered aspect ratio: 1:1
 - Intrinsic size: 128 × 128 px
 - Intrinsic aspect ratio: 1:1
 - File size: 3.2 kB
 - Current source: https://discords.com/_next/image?url=https%3A%2F%2Fcdn.discordapp.com%2Femojis%2F812334984928296986.png%3Fv%3D1&w=128&q=75

Information of an emoji that we can extract:

1. Image source (URL).
2. Title.

Data Collection and Preparation

❖ Analyze website

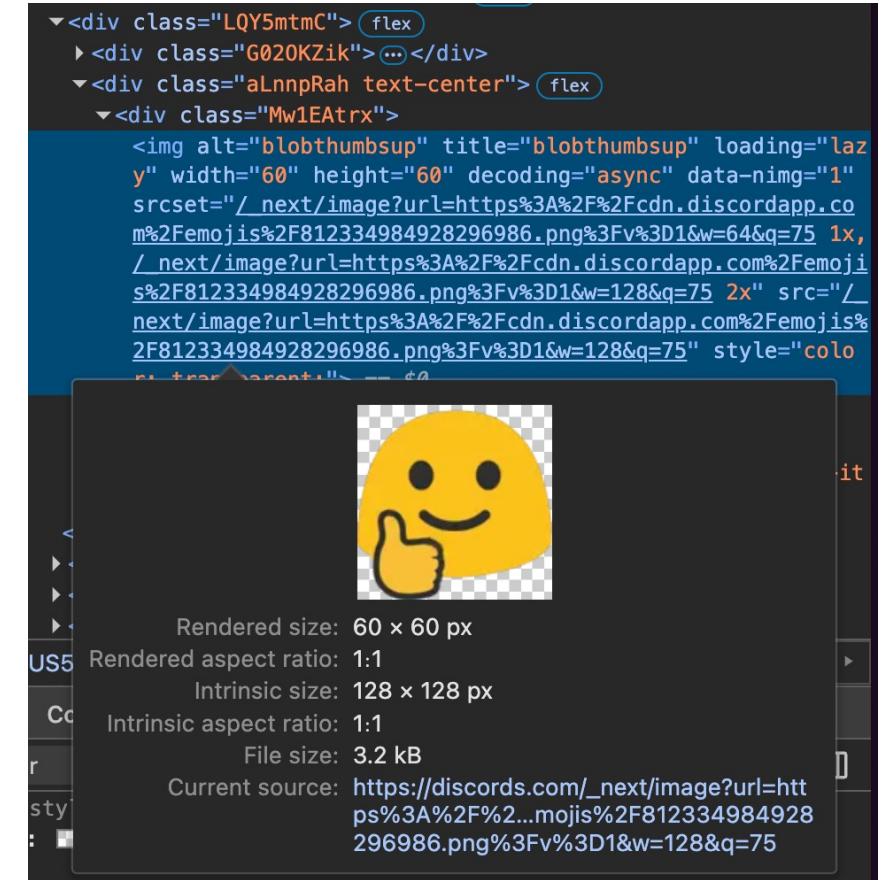


<https://discords.com/emoji-list/tag/Blobs?page=1>

Data Collection and Preparation

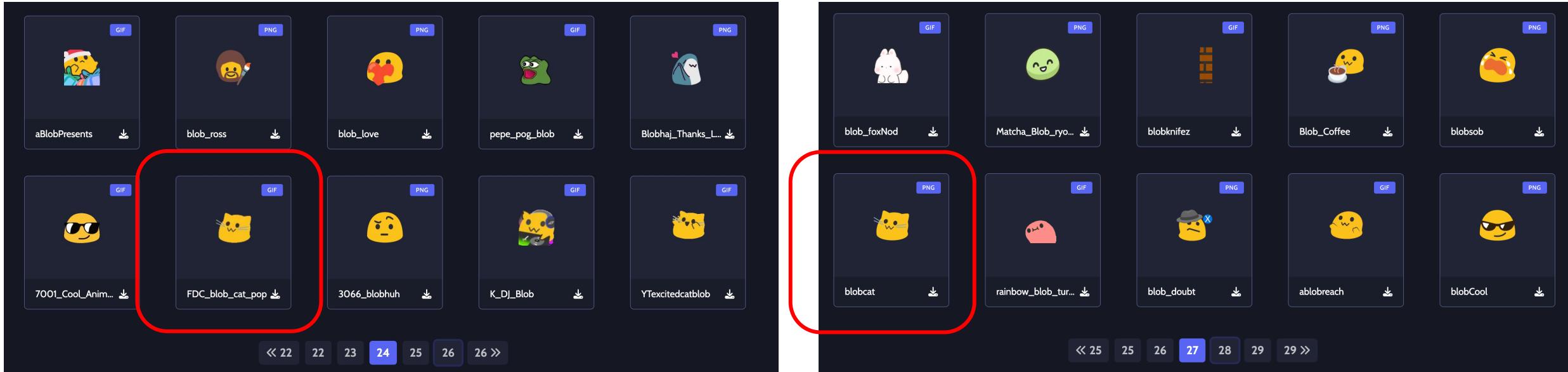
❖ Coding: Extract Image URLs

```
1 def get_image_links_from_page(page_url, driver):
2     driver.get(page_url)
3     try:
4         container = wait.until(EC.presence_of_element_located(
5             (By.CSS_SELECTOR, "div.FS5UE28h.container")
6         ))
7         image_items = wait.until(EC.presence_of_all_elements_located(
8             (By.CSS_SELECTOR, "div.LQY5mtmC div.aLnnpRah.text-center"))
9         )
10
11     image_links = []
12     for img_elem in image_items:
13         img_div = img_elem.find_element(By.CSS_SELECTOR, "div.Mw1EAttrx img, img")
14
15         img_url = img_div.get_attribute("src")
16         img_title = img_div.get_attribute("title")
17         if img_url:
18             image_links.append((img_url, img_title))
19
20     return image_links
21 except Exception as e:
22     print(f"Error while trying to extract images: {e}")
23     return []
```



Data Collection and Preparation

❖ Coding: Check image duplication



On this website, there might be the case where some emojis are duplicated.

Data Collection and Preparation

❖ Coding: Check image duplication

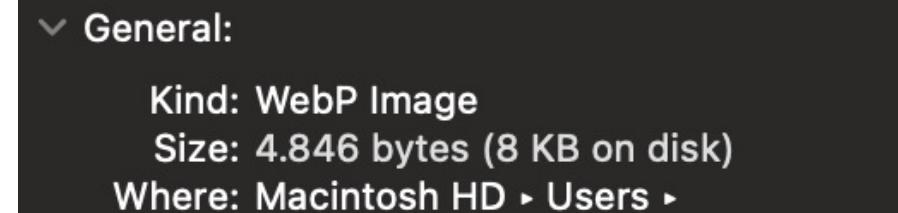
```
25 def hash_image_content(url):
26     try:
27         response = requests.get(url, stream=True)
28         if response.status_code == 200:
29             return hashlib.md5(response.content).hexdigest()
30         else:
31             print(f"Error downloading image from {url}; status: {response.status_code}")
32             return None
33     except requests.exceptions.RequestException as e:
34         print(f"Error with the image download for {url}: {e}")
35         return None
```

For every image content, convert it into a hexadecimal string that will be used as a **unique identifier** to check for duplication.

Data Collection and Preparation

❖ Coding: Convert image format

```
37 def convert_webp_to_jpg(webp_data):
38     try:
39         img = Image.open(BytesIO(webp_data))
40         if img.format == 'WEBP':
41             if img.mode == 'RGBA':
42                 img = img.convert('RGB')
43             buffer = BytesIO()
44             img.save(buffer, format="JPEG")
45             return buffer.getvalue()
46         else:
47             return webp_data
48     except Exception as e:
49         print(f"Error converting WebP to JPG: {e}")
50         return webp_data
```



Data Collection and Preparation

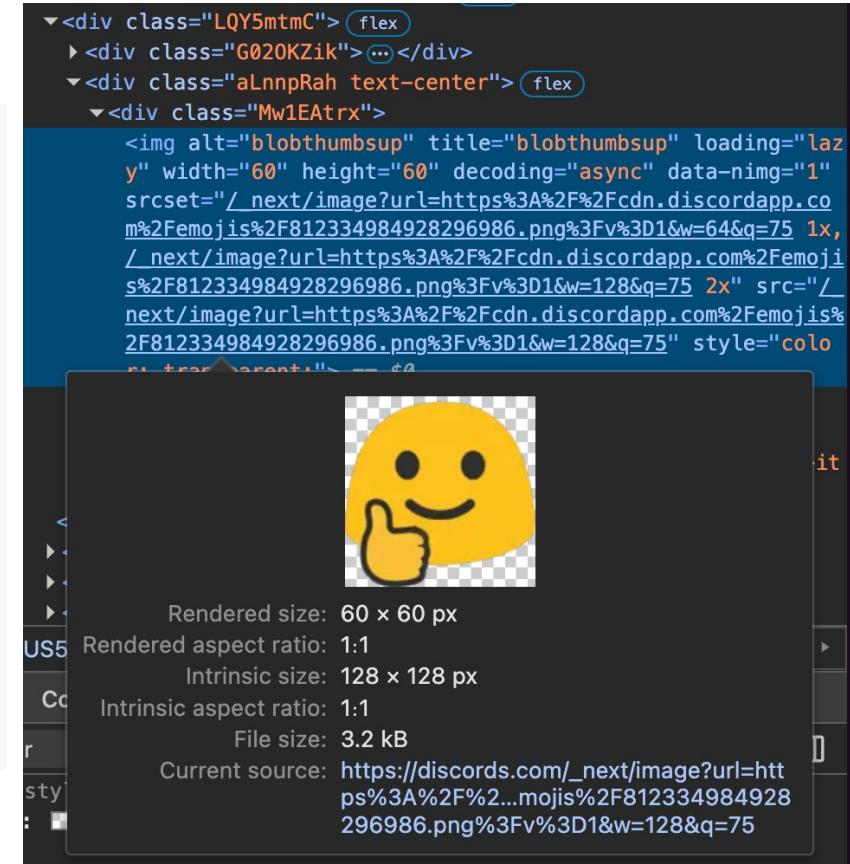
❖ Coding: Download image function

```
52 def download_image(img_url, img_name, folder_path):
53     try:
54         response = requests.get(img_url, stream=True)
55         if response.status_code == 200:
56             img_path = os.path.join(folder_path, f"{img_name}")
57             img_data = convert_webp_to_jpg(response.content)
58             with open(img_path, "wb") as f:
59                 f.write(img_data)
60         else:
61             print(f"Error downloading image from {img_url}; status: {response.status_code}")
62     except requests.exceptions.RequestException as e:
63         print(f"Error with the image download for {img_url}: {e}")
```

Data Collection and Preparation

❖ Coding: Process image function

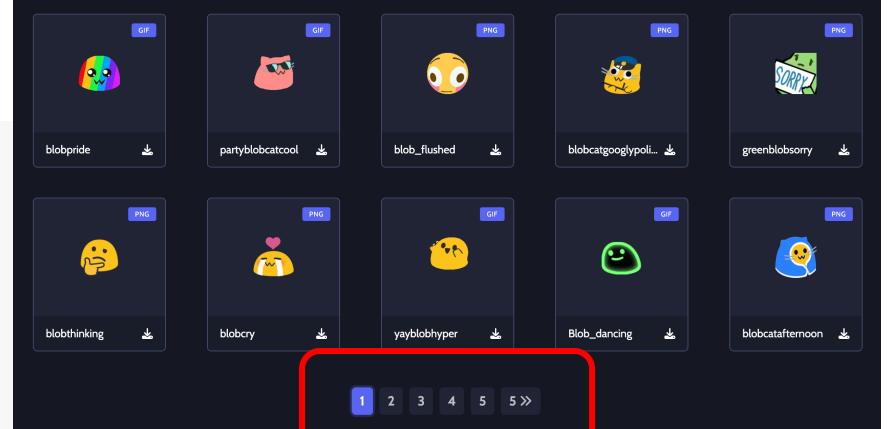
```
65 def process_image_page(image_url, img_title, folder_path, idx, tag, seen_hashes):
66     img_hash = hash_image_content(image_url)
67     if img_hash and img_hash not in seen_hashes:
68         seen_hashes.add(img_hash)
69         new_file_name = f"{tag}_{idx:07d}.jpg"
70         download_image(image_url, new_file_name, folder_path)
71         metadata = {
72             "file_name": new_file_name,
73             "image_url": image_url,
74             "image_title": img_title,
75             "tag": tag
76         }
77         return metadata
78     else:
79         return None
```



Data Collection and Preparation

❖ Coding: Process image function

```
81 def loop_over_pages(base_url, tags, total_pages, driver, folder_path):  
82     os.makedirs(folder_path, exist_ok=True)  
83     all_metadata = []  
84     seen_hashes = set()  
85  
86     for tag in tags:  
87         all_images = []  
88  
89         for page in tqdm(range(1, total_pages + 1), desc=f"Extracting Images for {tag}", unit="page"):  
90             page_url = f"{base_url}/emoji-list/tag/{tag}?page={page}"  
91             images = get_image_links_from_page(page_url, driver)  
92             all_images.extend(images)  
93  
94             time.sleep(1)  
95  
96             metadata_list = []  
97             for idx, (img_url, img_title) in enumerate(all_images, start=1):  
98                 metadata = process_image_page(img_url, img_title, folder_path, idx, tag, seen_hashes)  
99                 if metadata:  
100                     metadata_list.append(metadata)  
101  
102             all_metadata.extend(metadata_list)  
103  
104     return all_metadata
```



```
106 def save_metadata(metadata_list, metadata_file):  
107     df = pd.DataFrame(metadata_list)  
108     df.to_csv(metadata_file,  
109                index=False,  
110                encoding="utf-8")
```

Data Collection and Preparation

❖ Coding: Start crawling

```
1 os.makedirs("crawled_data", exist_ok=True)
2 folder_path = os.path.join("crawled_data", "images")
3 metadata_file = os.path.join("crawled_data", "metadata.csv")
4
5 base_url = "https://discords.com"
6 tags = ["Panda"]
7 total_pages = 1000
8
9 metadata_list = loop_over_pages(base_url, tags, total_pages, driver, folder_path)
10 save_metadata(metadata_list, metadata_file)
11
12 print("Start downloading images...")
13 with tqdm(total=len(metadata_list), desc="Downloading Images", unit="image") as pbar:
14     for metadata in metadata_list:
15         img_url = metadata['image_url']
16         file_name = metadata['file_name']
17         download_image(img_url, file_name, folder_path)
18         pbar.update(1)
19
20 print("Download images completed.")
21
22 total_crawled_images = len(os.listdir(folder_path))
23 print(f"Total crawled images: {total_crawled_images}.")
24
25 driver.quit()
```

```
#11 0x592a2065bce7 <unknown>
#12 0x592a2065b9e4 <unknown>
#13 0x592a2066013a <unknown>
#14 0x592a2065c5b9 <unknown>
#15 0x592a20641e00 <unknown>
#16 0x592a206735d2 <unknown>
#17 0x592a20673778 <unknown>
#18 0x592a2068ba1f <unknown>
#19 0x7ae3c963fac3 <unknown>
#20 0x7ae3c96d1850 <unknown>
```

```
Extracting Images for Panda: 100%|██████████| 3942/3942
Error downloading image from https://discords.com
Start downloading images...
Downloading Images: 100%|██████████| 3942/3942
Download images completed.
Total crawled images: 3942.
```

Data Collection and Preparation

❖ Crawled emojis

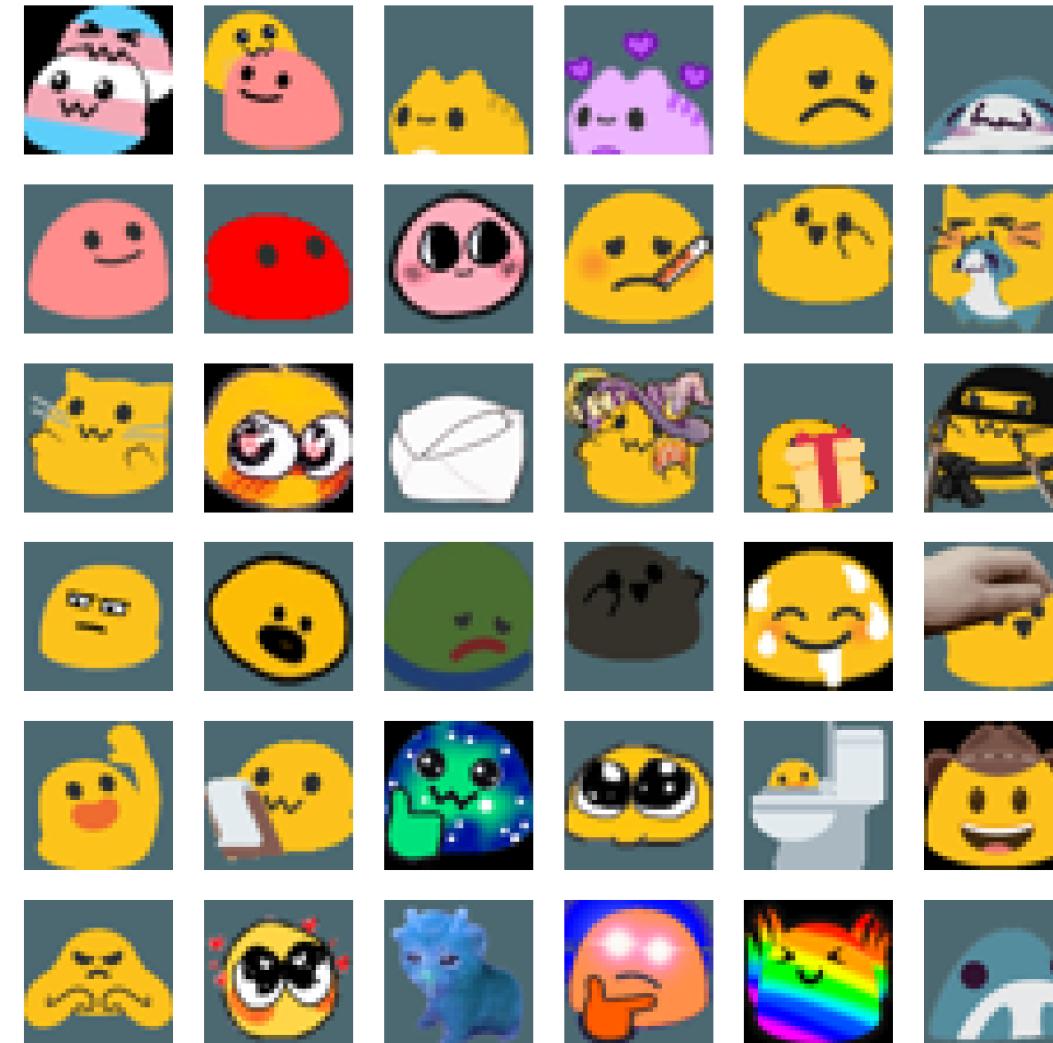
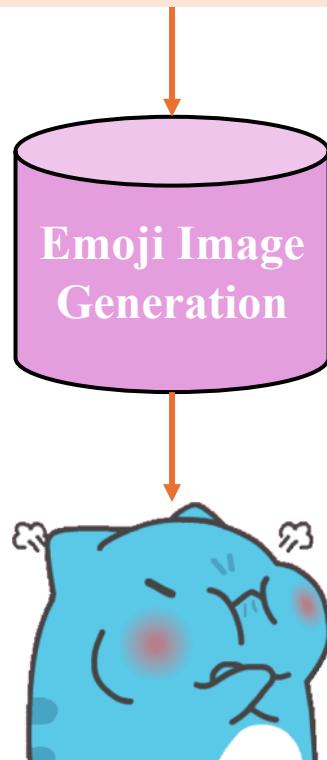


Figure: A bunch of emojis crawled using our code.

Data Collection and Preparation

❖ Crawled emojis

“Blue cat emoji with a round head, blushing cheeks, puffed face, angry expression, tiny clenched fists raised beside the head, steam puffs above both ears, no background.”



Input

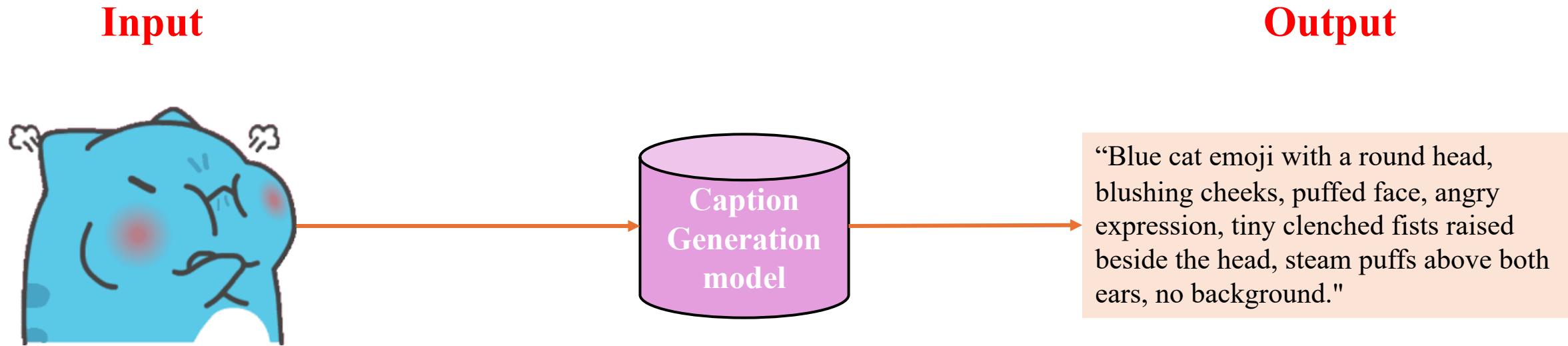
Output



Currently, there is no complete prompt/description for each image of our dataset.

Data Collection and Preparation

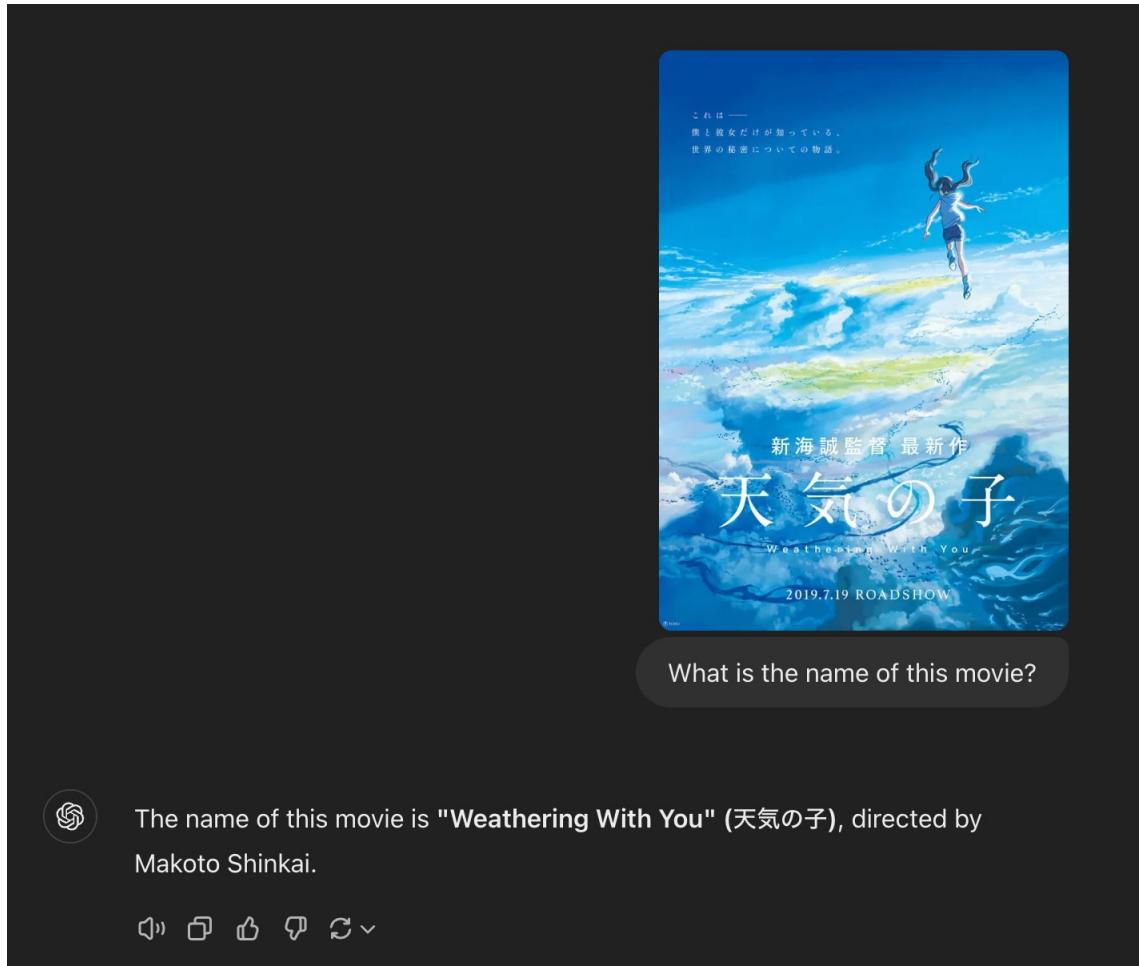
❖ Caption for each emoji



Thus, to complete the dataset where each sample consists of an image and a prompt, a model to create pseudo label for image is required.

Data Collection and Preparation

❖ Idea: Use LVLM?



Large Vision-Language Models (LVLMs) are models that integrate large language models (LLMs) with pre-trained vision encoders, enabling them to process image inputs, understand various queries, and perform reasoning tasks.

Figure: ChatGPT, one of the most famous application that utilizing LVLMs.

Data Collection and Preparation

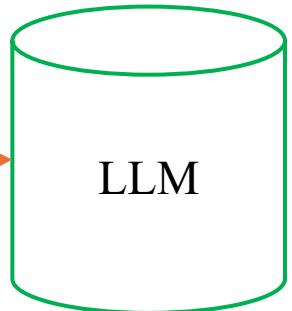
❖ Idea: Use LVLM?

Image:



Question: What is the feeling of the cat in this image?

?



LLM only accepts text, how can an image be feeded into?

Figure: Illustration of a pair of (image, text) input.

Data Collection and Preparation

❖ LVLMs Architecture

Output:

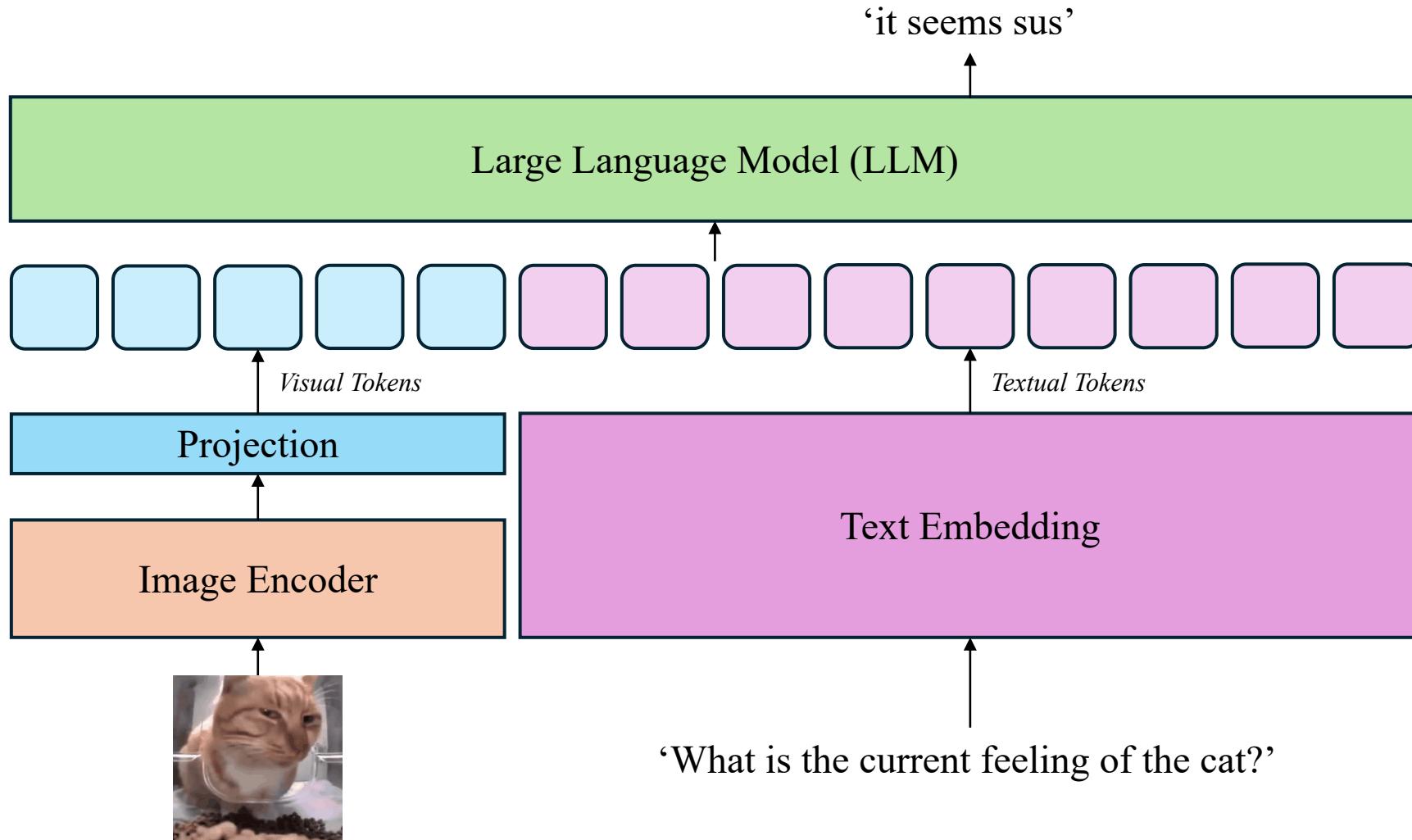


Figure: A typical architecture of LVLMs.

Input:

Data Collection and Preparation

❖ Capability of VLMs

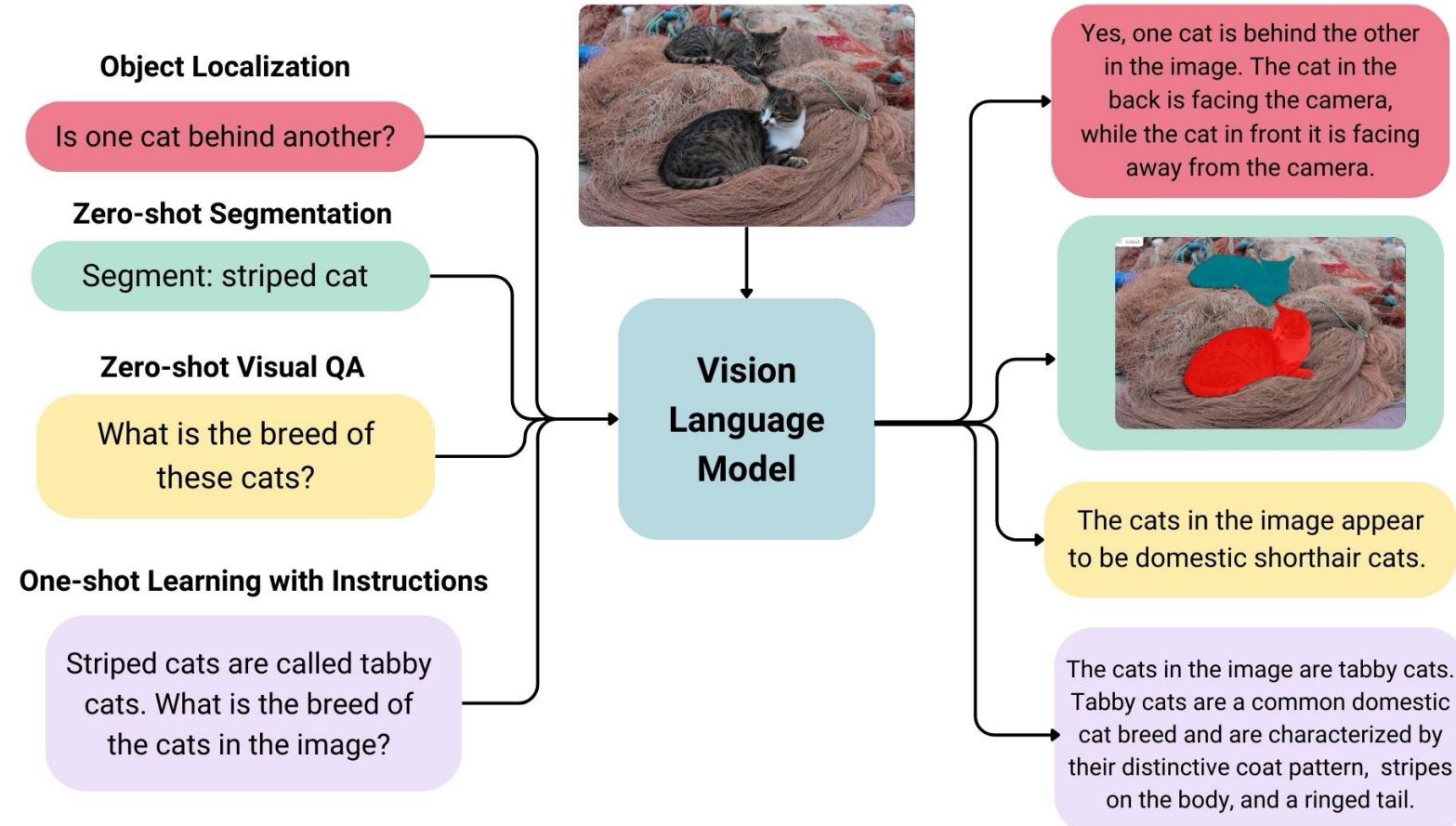
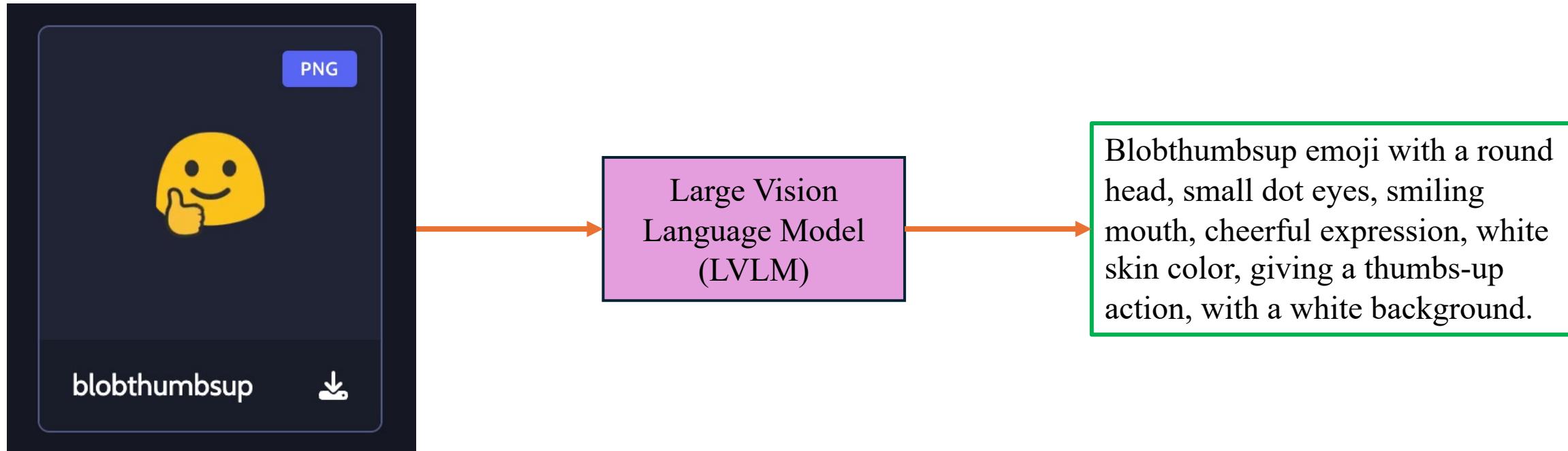


Figure: Illustration of what a Large Vision-Language Model (LVLM) can do nowadays.

Data Collection and Preparation

❖ Generate caption for emojis



Data Collection and Preparation

❖ Choosing LVLMs



Qwen-VL

Qwen-VL: A large-scale vision-language model that integrates image and text understanding for multi-modal tasks, enabling advanced capabilities such as image captioning, visual question answering, and cross-modal reasoning.

It leverages transformer-based architectures and extensive training to generate highly coherent outputs and accurately link textual information with visual content.



Data Collection and Preparation

❖ Choosing LLMs

Qwen/Qwen2.5-VL-3B-Instruct □ like 319 Follow Qwen 24.8k

Image-Text-to-Text Transformers Safetensors English qwen2_5_vl multimodal conversational text-generation-inference arxiv:2309.00071

arxiv:2409.12191 arxiv:2308.12966

Model card Files and versions Community 24 Edit model card

Train Deploy Use this model

Qwen2.5-VL-3B-Instruct

Qwen Chat

Introduction

In the past five months since Qwen2-VL's release, numerous developers have built new models on the Qwen2-VL vision-language models, providing us with valuable feedback. During this period, we focused on building more useful vision-language models. Today, we are excited to introduce the latest addition to the Qwen family: Qwen2.5-VL.

Key Enhancements:

- Understand things visually: Qwen2.5-VL is not only proficient in recognizing common objects such as flowers, birds, fish, and insects, but it is highly capable of analyzing texts, charts, icons, graphics, and layouts within images.

Downloads last month 1,204,016

Safetensors Model size 3.75B params Tensor type BF16

Inference Providers

Image-Text-to-Text

This model isn't deployed by any Inference Provider. Ask for provider support

Model tree for Qwen/Qwen2.5-VL-3B-Instruct

- Adapters 12 models
- Finetunes 106 models
- Quantizations 25 models

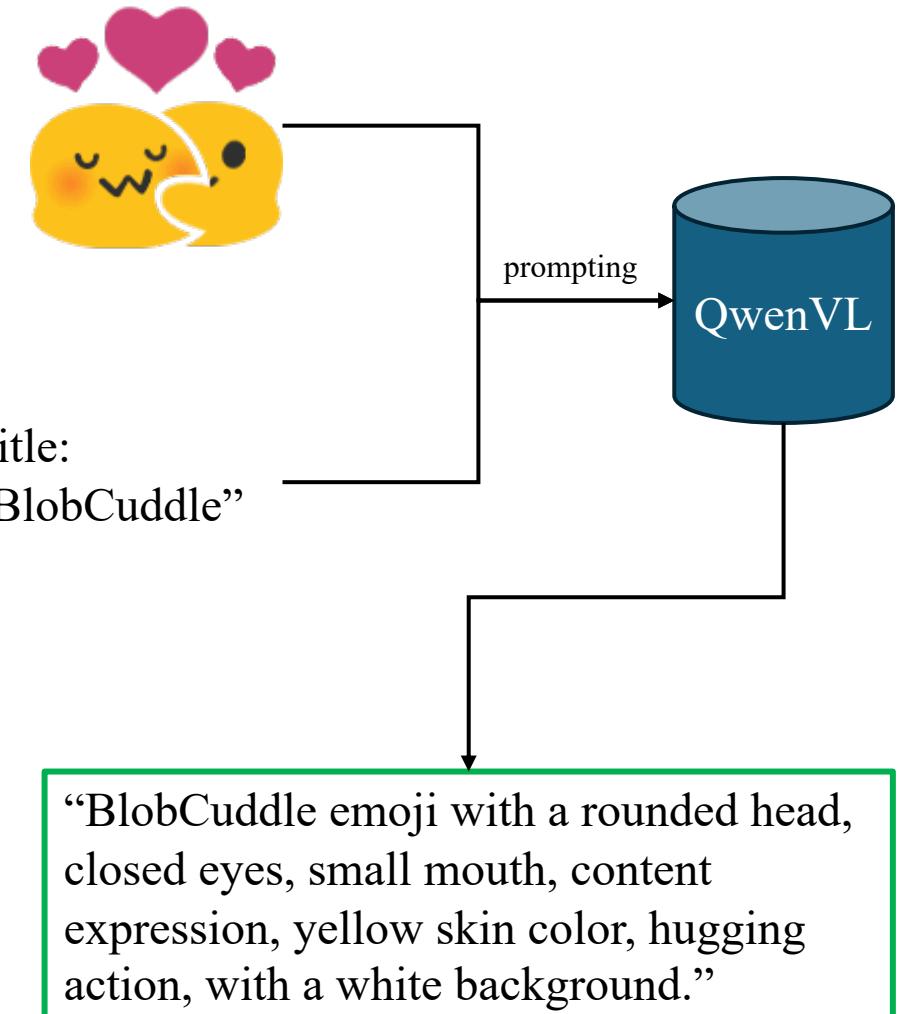
Data Collection and Preparation

❖ Prompting QwenVL for emoji captioning

Instruction: Given the emoji title and the corresponding image, generate a natural language description of the emoji's key features. The description should include: 1) the head shape, 2) eye characteristics, 3) mouth characteristics, 4) facial expression, 5) skin color, 6) any action (if present), and 7) background color. The description should be concise and structured like this: '{emoji_title} emoji with a {head shape}, {eye description} eyes, {mouth description} mouth, {expression description} expression, {skin color} skin color, {action description}, with a {background color} background.'

Example: For an emoji titled 'Pepe the Frog', the description should be like: 'Pepe the frog emoji with a round head, big eyes, smiling mouth, happy expression, green skin color, no action, with a green background.'

Input: '{emoji_title}'t



Data Collection and Preparation

❖ Coding

```
1 import os
2 import pandas as pd
3 import torch
4 from transformers import Qwen2_5_VLForConditionalGeneration
5 from transformers import AutoProcessor
6 from PIL import Image
7 from qwen_vl_utils import process_vision_info
8 from tqdm import tqdm
```

```
1 model_id = "Qwen/Qwen2.5-VL-3B-Instruct"
2
3 model = Qwen2_5_VLForConditionalGeneration.from_pretrained(
4     model_id,
5     torch_dtype=torch.bfloat16,
6     device_map="auto"
7 )
8
9 processor = AutoProcessor.from_pretrained(model_id)
```



Qwen-VL

Data Collection and Preparation

❖ Coding: Create prompt function

```
1 def create_prompt(emoji_title):
2     instruction = (
3         "Instruction: Given the emoji title and the corresponding image, generate a natural language description "
4         "of the emoji's key features. The description should include: 1) the head shape, 2) eye characteristics, "
5         "3) mouth characteristics, 4) facial expression, 5) skin color, 6) any action (if present), and 7) background color. "
6         "The description should be concise and structured like this: "
7         "'{emoji_title} emoji with a {head shape}, {eye description} eyes, {mouth description} mouth, "
8         "{expression description} expression, {skin color} skin color, {action description}, "
9         "with a {background color} background.'"
10    )
11
12    example = (
13        "Example: For an emoji titled 'Pepe the Frog', the description should be like: "
14        "'Pepe the frog emoji with a round head, big eyes, smiling mouth, happy expression, "
15        "green skin color, no action, with a green background.'"
16    )
17
18    prompt = f"{instruction}\n{example}\nInput: '{emoji_title}'"
19    return prompt
```

Data Collection and Preparation

❖ Coding: Setup input

```
for _, row in tqdm(df.iterrows(), total=len(df), desc=f"Processing {subfolder}"):
    image_path = os.path.join(images_folder, row["file_name"])
    prompt = create_prompt(row["image_title"])

    image = Image.open(image_path).convert("RGBA")

    messages = [
        {
            "role": "user",
            "content": [
                {"type": "image", "image": image},
                {"type": "text", "text": prompt},
            ],
        }
    ]

    text = processor.apply_chat_template(
        messages, tokenize=False, add_generation_prompt=True
    )
    image_inputs, video_inputs = process_vision_info(messages)
```

Data Collection and Preparation

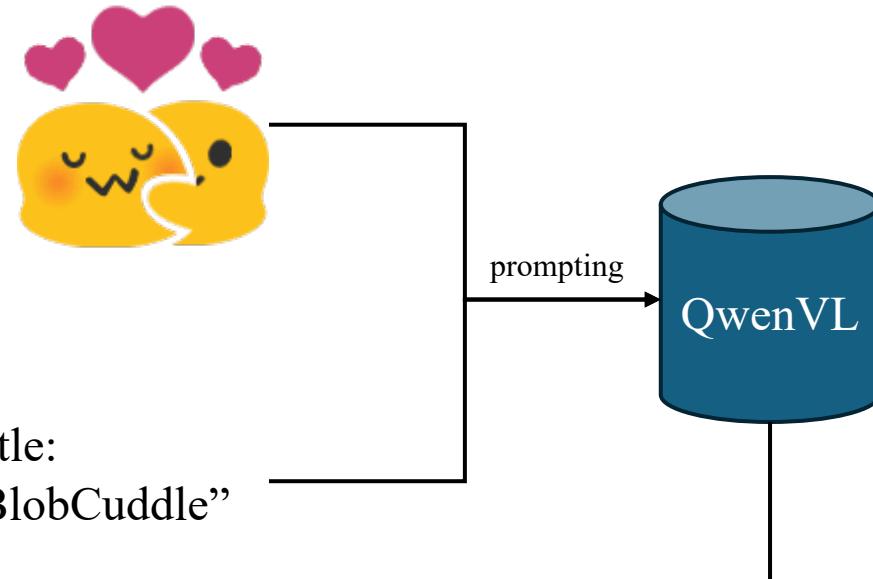
❖ Coding: Generate caption

```
inputs = processor(  
    text=[text],  
    images=image_inputs,  
    videos=video_inputs,  
    padding=True,  
    return_tensors="pt"  
).to("cuda")  
  
generated_ids = model.generate(**inputs, max_new_tokens=128)  
generated_ids_trimmed = [  
    out_ids[len(in_ids):] for in_ids, out_ids in zip(inputs.input_ids, generated_ids)  
]  
result_text = processor.batch_decode(  
    generated_ids_trimmed, skip_special_tokens=True, clean_up_tokenization_spaces=False  
)  
  
df.loc[_, "prompt"] = result_text
```

Data Collection and Preparation

❖ Results

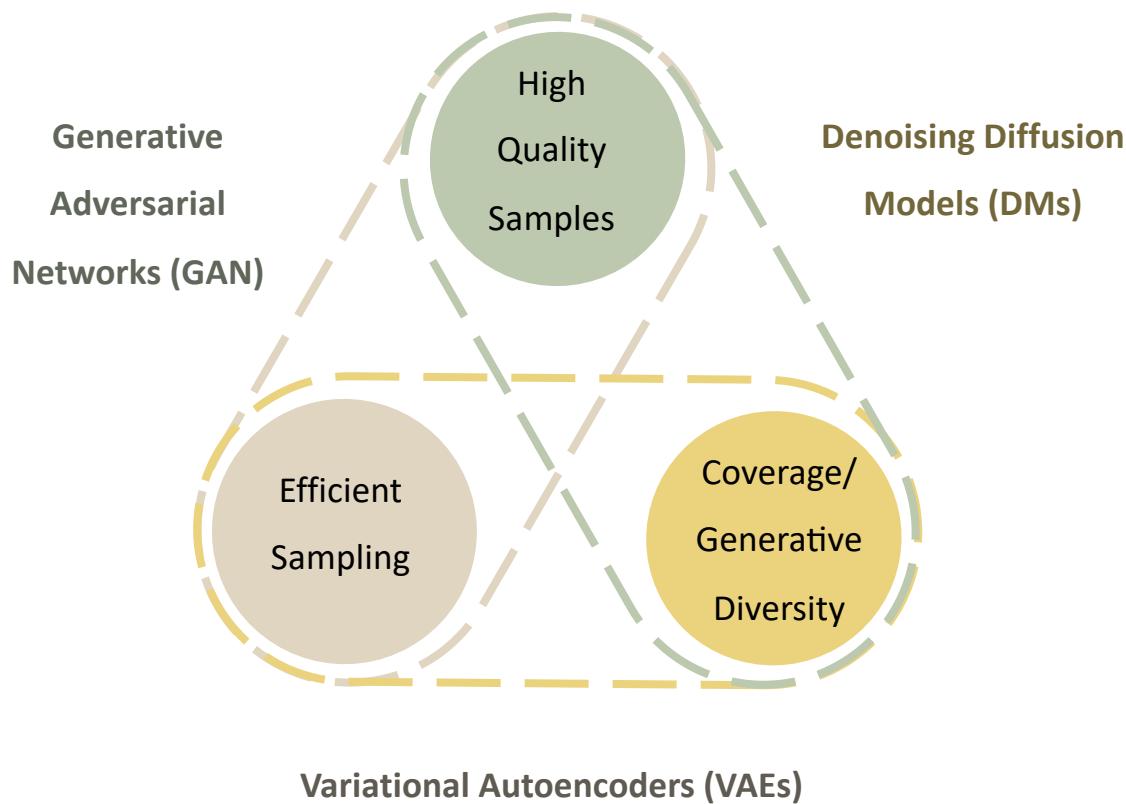
image_title	tag	prompt
bloboof	Blobs	bloboof emoji with a square head, small eyes, open mouth, surprised expression, white skin color, no action, with a white background.
blobhyperthink	Blobs	blobhyperthink emoji with a round head, large circular eyes, wide open mouth, surprised expression, orange skin color, no action, with a white background.
blob_minecraft_totem_of_undying	Blobs	blob_minecraft_totem_of_undying emoji with a blob-like head, large eyes, wide open mouth, surprised expression, yellow skin color, no action, with a white background.
BlobCuddle	Blobs	BlobCuddle emoji with a rounded head, closed eyes, small mouth, content expression, yellow skin color, hugging action, with a white background.
nyan_party_blob	Blobs	Nyan Party Blob emoji with a rounded head, small black eyes, squiggly mouth, surprised expression, pink skin color, no action, with a white background.
purple_blob	Blobs	Purple Blob emoji with a round head, small black eyes, smiling mouth, happy expression, purple skin color, waving action, with a white background.



Training SD from scratch

Training SD from scratch

❖ Introduction: Generative models

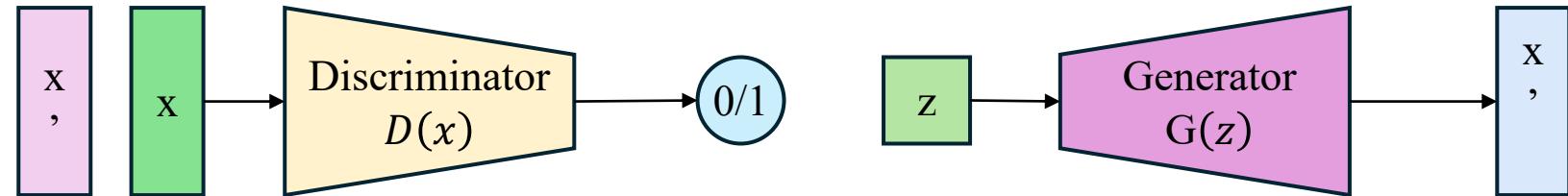


FEATURES	GAN	VAE	DIFFUSION MODEL
IMAGE QUALITY	High-quality, sharp images	Often blurry, less realistic	Very high and realistic
SAMPLE DIVERSITY	Lower	Moderate diversity	High diversity
SAMPLING SPEED	Fast	Fast	Slow

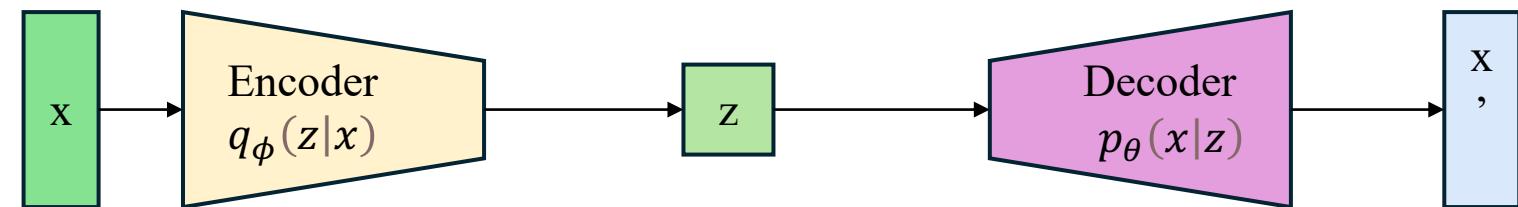
Training SD from scratch

❖ Introduction: Generative models

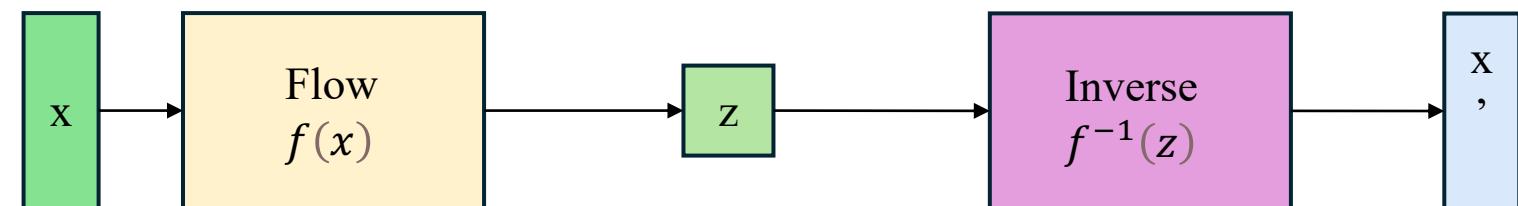
GAN: Adversarial training.



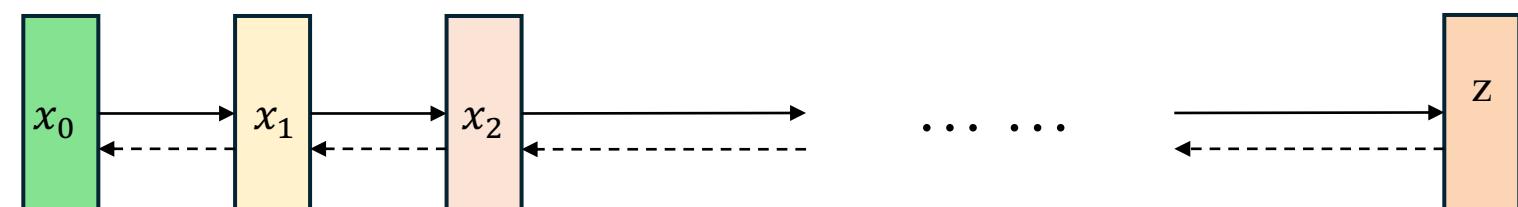
VAE: maximize variational lower bound.



Flow-based models:
Invertible transform of distributions.



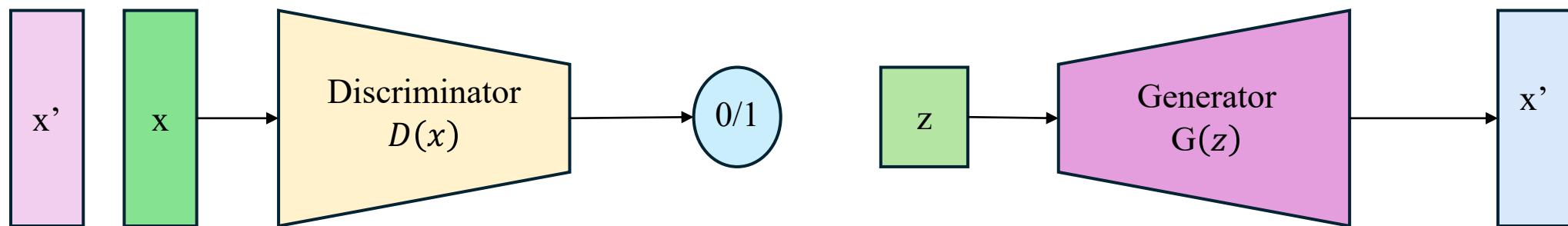
Diffusion models: Gradually add Gaussian noise and then reverse.



Training SD from scratch

❖ Introduction: GANs

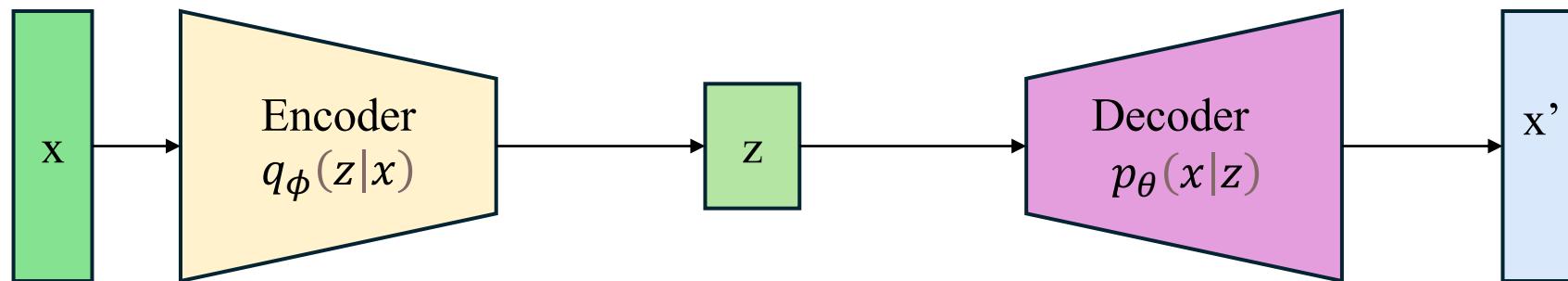
GAN: Adversarial training



Training SD from scratch

❖ Introduction: VAE

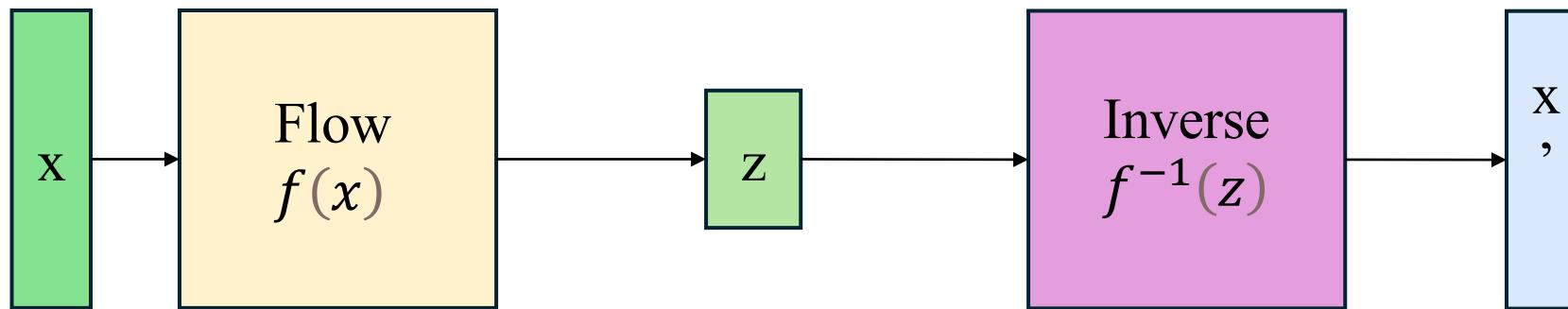
VAE: maximize variational
lower bound.



Training SD from scratch

❖ Introduction: Flow Matching

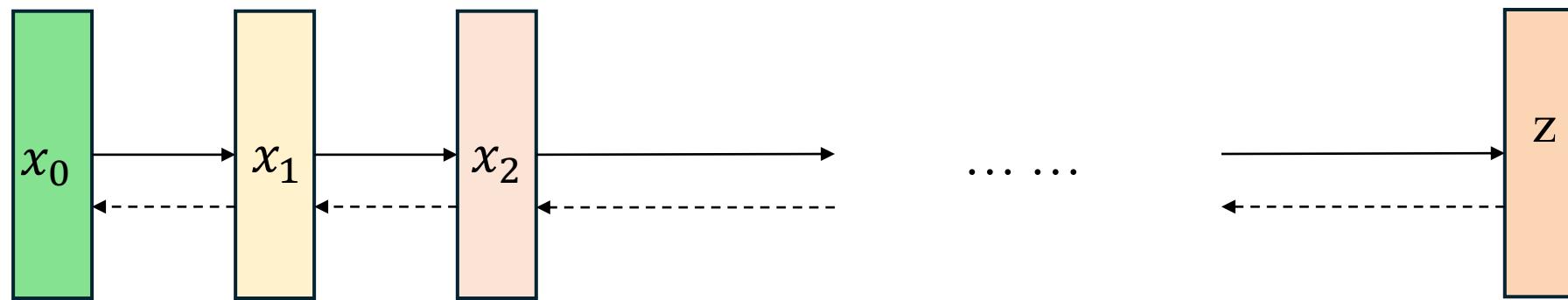
Flow-based models:
Invertible transform of distributions.



Training SD from scratch

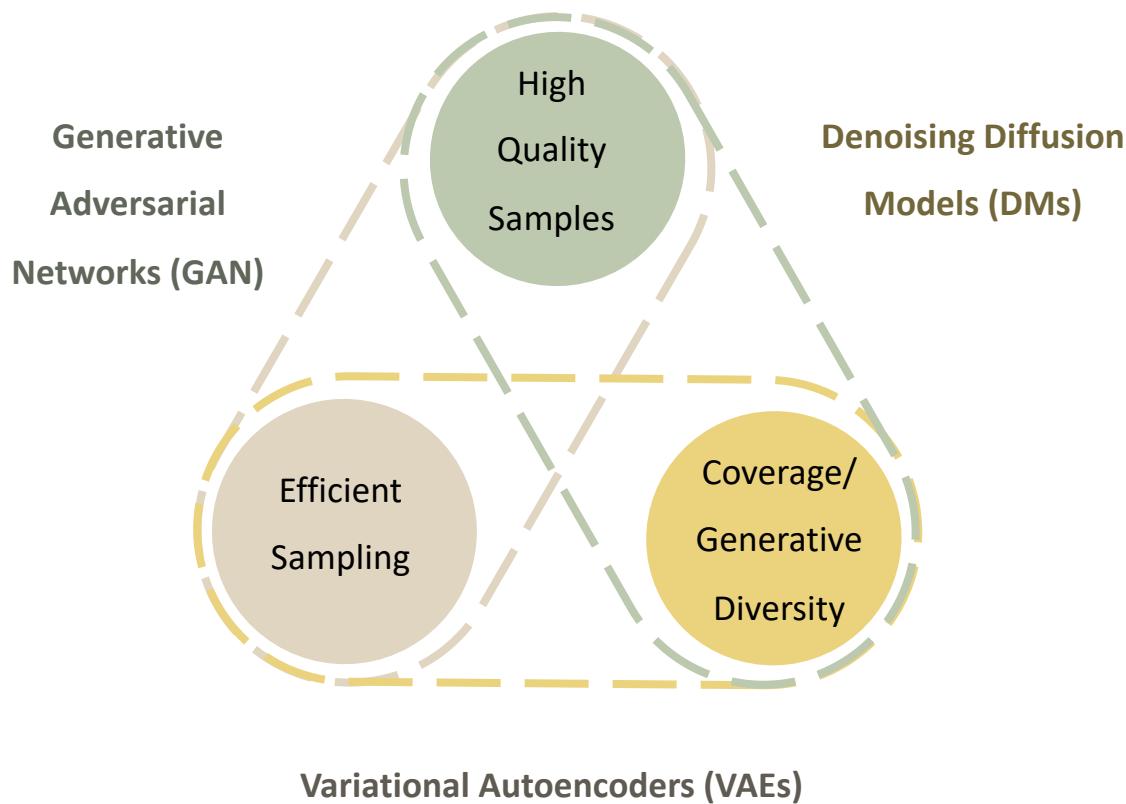
❖ Introduction: Diffusion models

Diffusion models: Gradually add Gaussian noise and then reverse.



Training SD from scratch

❖ Introduction: Generative models

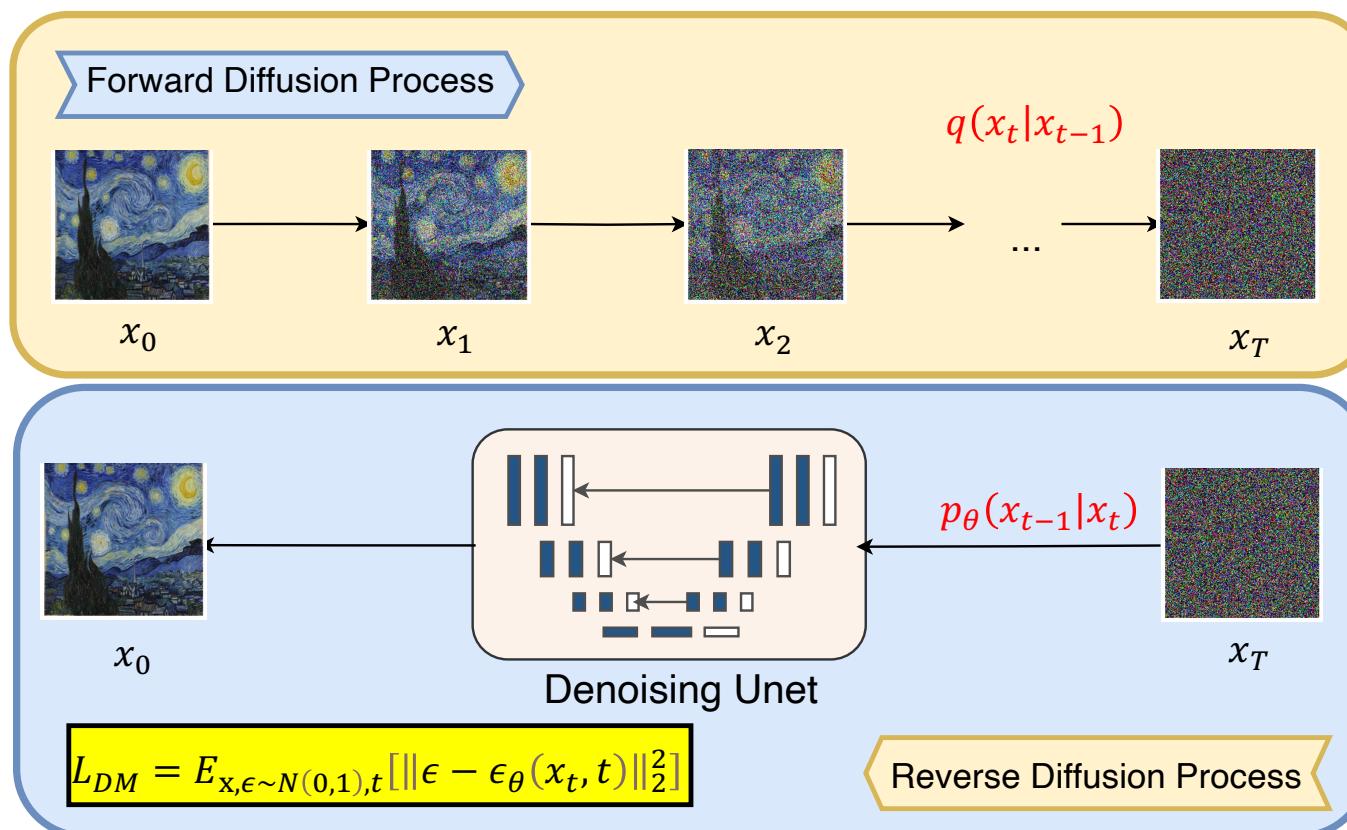


FEATURES	GAN	VAE	DIFFUSION MODEL
IMAGE QUALITY	High-quality, sharp images	Often blurry, less realistic	Very high and realistic
SAMPLE DIVERSITY	Lower	Moderate diversity	High diversity
SAMPLING SPEED	Fast	Fast	Slow

Why?

Training SD from scratch

❖ Introduction: Pure Diffusion Models



Limitations:

- Requires **hundreds of sampling steps**.
- **Resource-intensive and very slow.**

This speed issue led to the development of
Stable Diffusion!

Training SD from scratch

❖ Introduction: An example of Stable Diffusion

An original card from the Rider-Waite Tarot deck (1909).

(Source: [Wikipedia](#))



A fantasy-inspired style, drawing from watercolor paintings and mythological elements. (Generated by AI - Stable Diffusion 3.5)



Stable diffusion model

Prompt

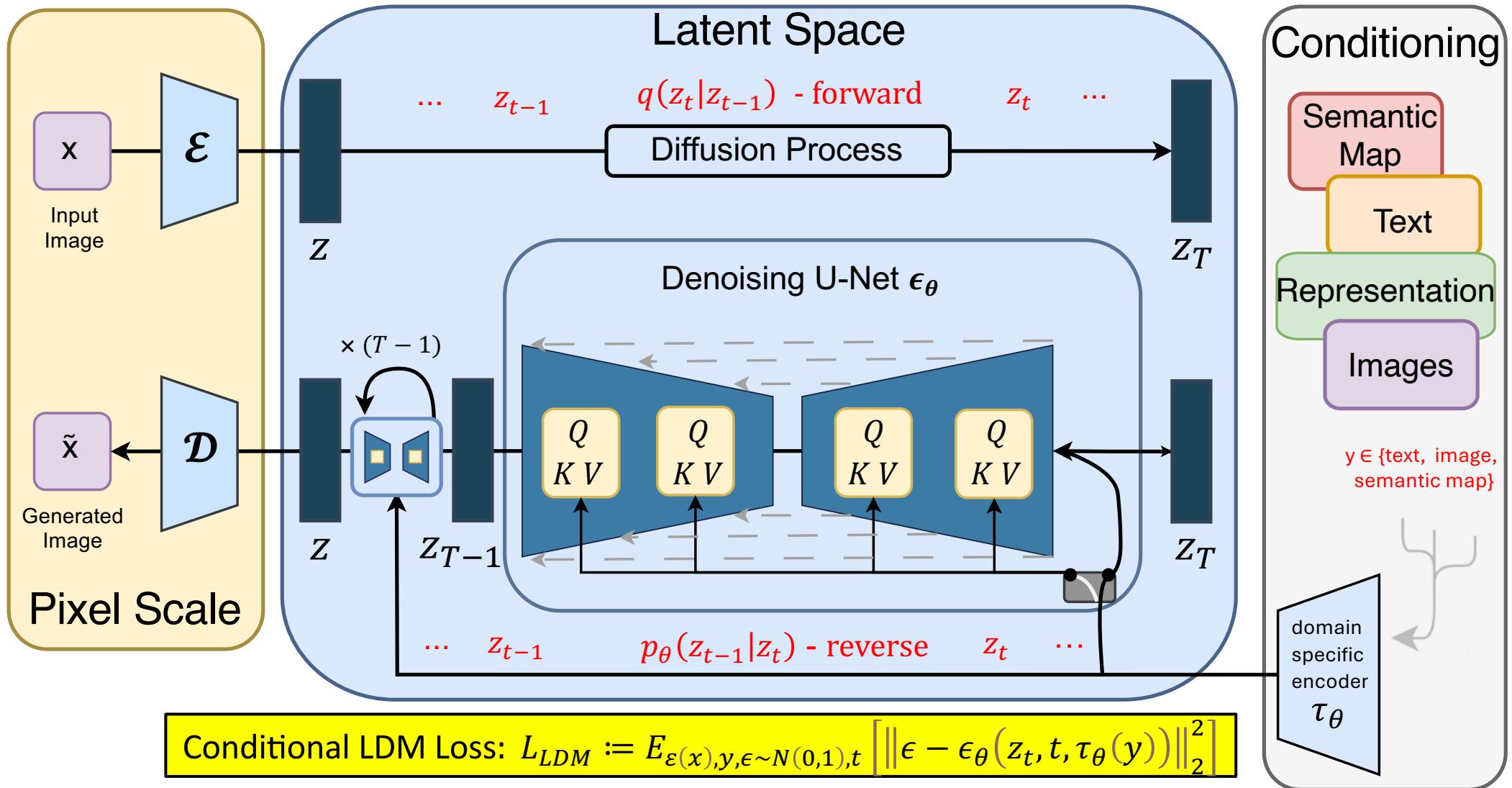
“An illustration of the Tarot card *The Empress*, preserving the composition and symbols of the Rider-Waite deck, but reimagined in a fantasy style inspired by watercolor art and mythical elements.”



The new illustrated version of the card adopts a fantasy and mythological art style.

(Generated by AI - Stable Diffusion 3.5)

❖ Stable Diffusion Models Architecture



denoising step



cross-attention



switch

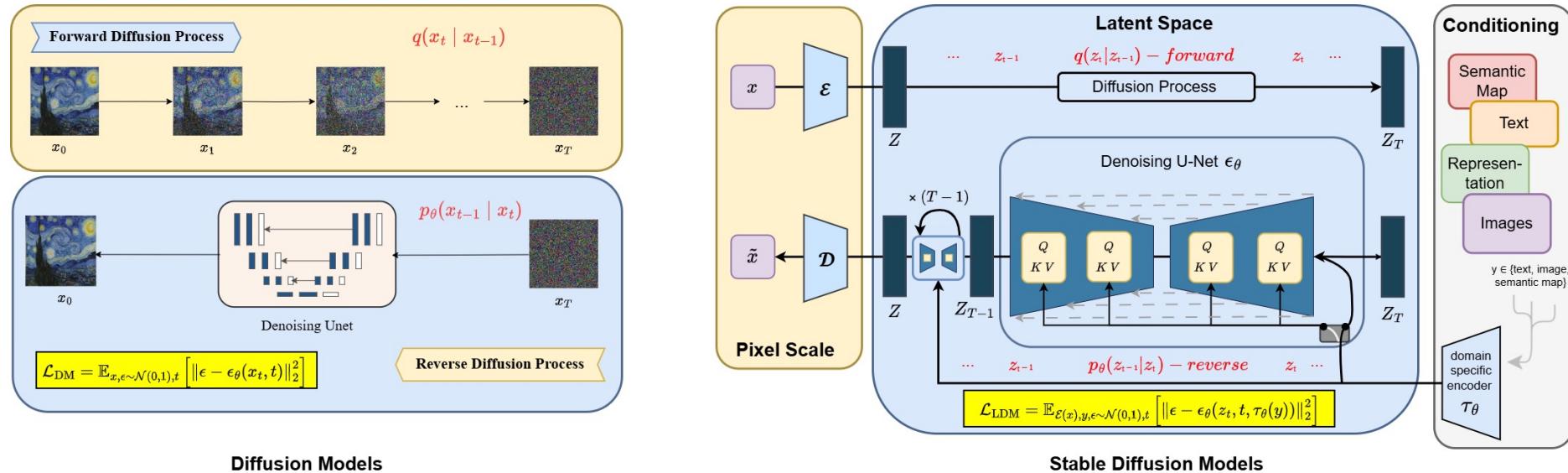


skip connection



concat

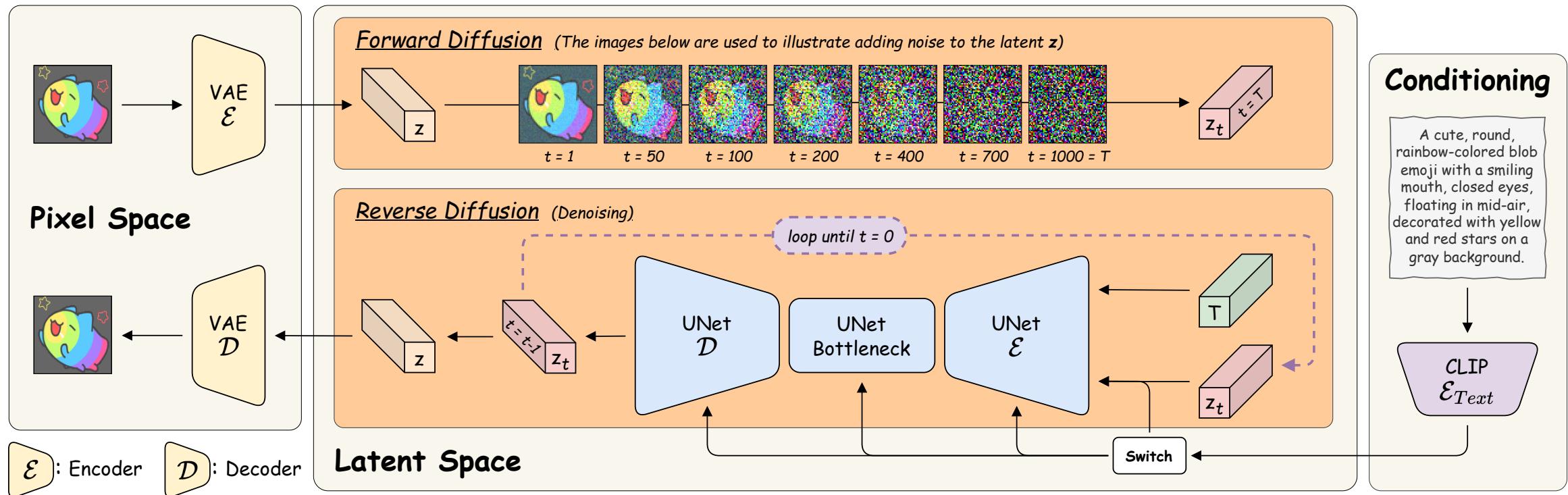
❖ Pure Diffusion Models vs Stable Diffusion Models



Features	Pure Diffusion Models	Stable Diffusion Models
Processing Space	Pixel space (e.g. $x \in R^{H \times W \times 3}$)	Latent space $z \in R^d$ (via VAE encoder)
Inference Speed	Slower	Faster
Architecture Highlight	Denoising U-Net only	Denoising U-Net with Cross-Attention for prompts

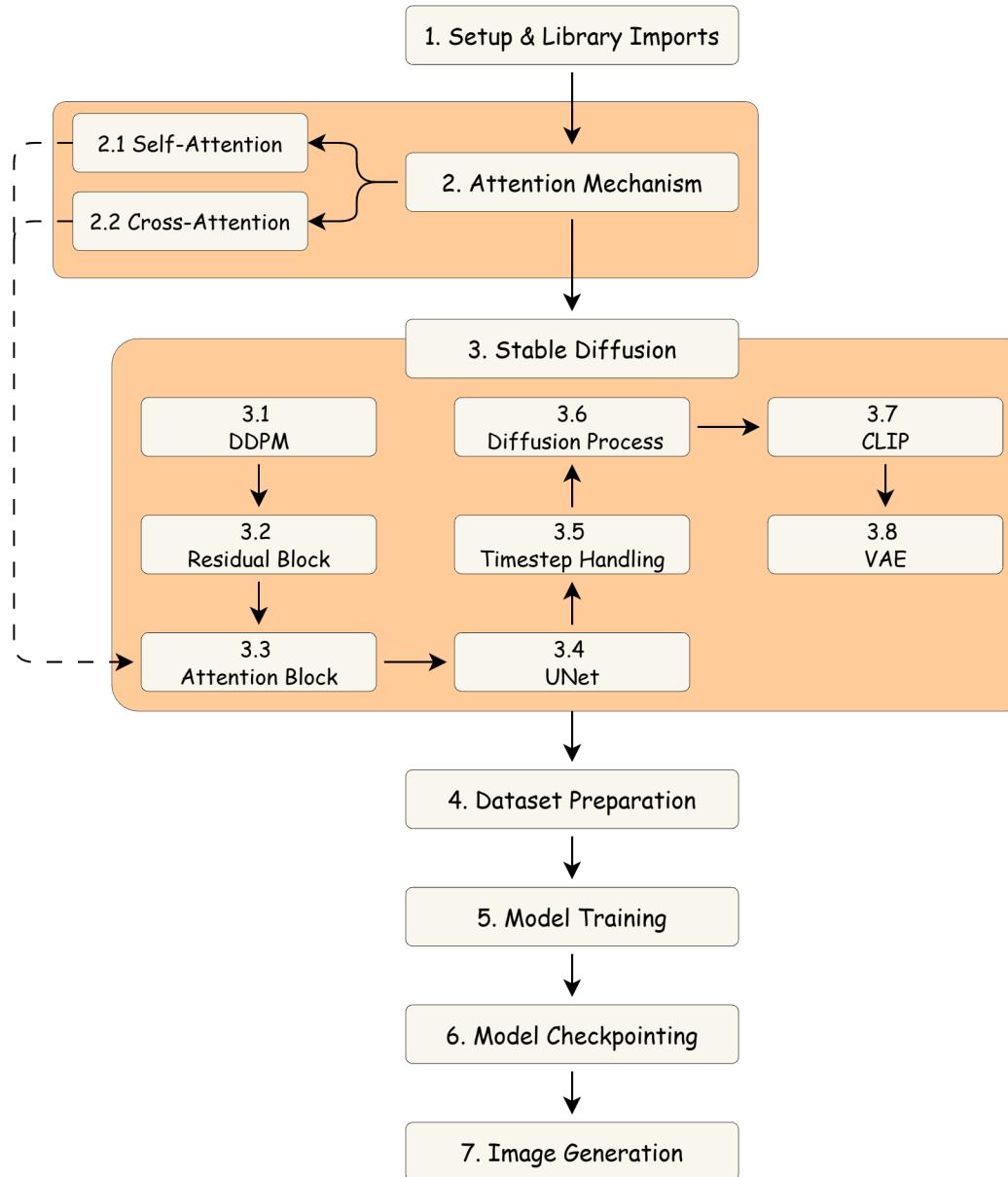
Training SD from scratch

❖ Overall architecture



Training SD from scratch

❖ Workflow



Training SD from scratch

❖ Setup & Library Imports



Dynamiteffusers



Transformers

```
1 import os
2 import math
3 import torch
4 import random
5 import numpy as np
6 import pandas as pd
7 import matplotlib.pyplot as plt
8
9 from torch import nn
10 from PIL import Image
11 from tqdm import tqdm
12 from torchvision import transforms
13 from diffusers import AutoencoderKL
14 from torch.nn import functional as F
15 from torch.amp import GradScaler, autocast
16 import torch.optim.lr_scheduler as lr_scheduler
17 from torch.utils.data import Dataset, DataLoader
18 from transformers import CLIPTokenizer, CLIPTextModel
```

Training SD from scratch

❖ Diffuser



D ffusers

Popular Tasks & Pipelines

Task	Pipeline	😊 Hub
Unconditional Image Generation	DDPM	google/ddpm-ema-church-256
Text-to-Image	Stable Diffusion Text-to-Image	stable-diffusion-v1-5/stable-diffusion-v1-5
Text-to-Image	unCLIP	kakaobrain/karlo-v1-alpha
Text-to-Image	DeepFloyd IF	DeepFloyd/IF-I-XL-v1.0
Text-to-Image	Kandinsky	kandinsky-community/kandinsky-2-2-decoder
Text-guided Image-to-Image	ControlNet	Illyasviel/sd-controlnet-canny
Text-guided Image-to-Image	InstructPix2Pix	timbrooks/instruct-pix2pix
Text-guided Image-to-Image	Stable Diffusion Image-to-Image	stable-diffusion-v1-5/stable-diffusion-v1-5
Text-guided Image Inpainting	Stable Diffusion Inpainting	runwayml/stable-diffusion-inpainting
Image Variation	Stable Diffusion Image Variation	lambdalabs/sd-image-variations-diffusers
Super Resolution	Stable Diffusion Upscale	stabilityai/stable-diffusion-x4-upscaler
Super Resolution	Stable Diffusion Latent Upscale	stabilityai/sd-x2-latent-upscaler



InstructPix2Pix

<https://huggingface.co/timbrooks/instruct-pix2pix>



x4 Upscaler

<https://huggingface.co/stabilityai/stable-diffusion-x4-upscaler>

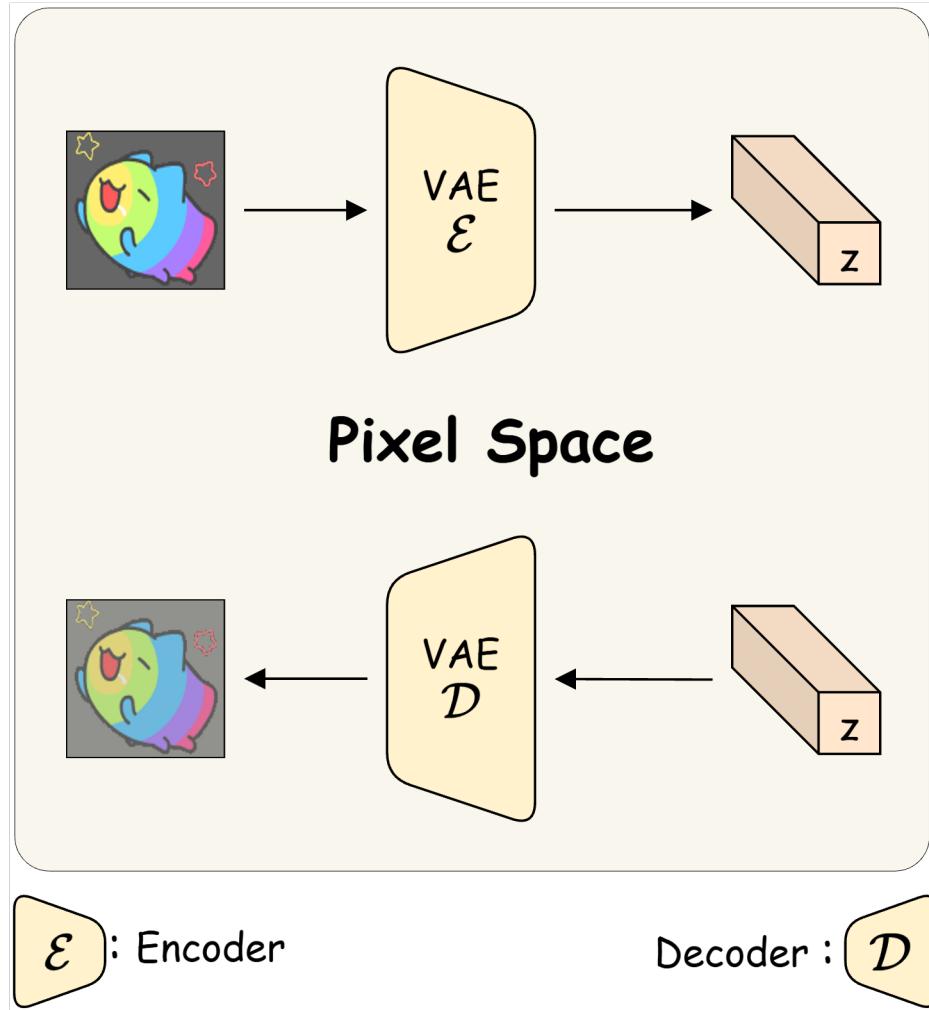


Stable Diffusion Image Variations

<https://huggingface.co/lambdalabs/sd-image-variations-diffusers>

Training SD from scratch

❖ VAE Review

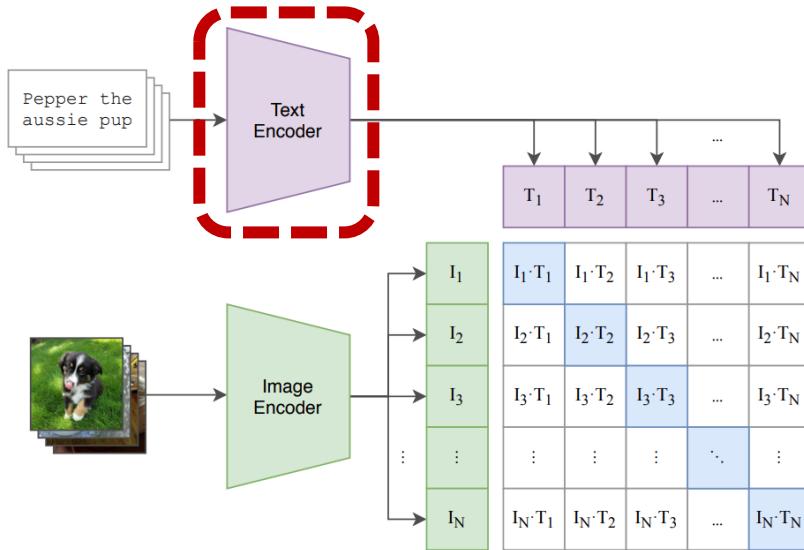


```
1 VAE_id = "stabilityai/sd-vae-ft-mse"  
2 vae = AutoencoderKL.from_pretrained(VAE_id)  
3 vae.requires_grad_(False)  
4 vae.eval()
```

Training SD from scratch

❖ CLIP Review

(1) Contrastive pre-training



Conditioning

A cute, round, rainbow-colored blob emoji with a smiling mouth, closed eyes, floating in mid-air, decorated with yellow and red stars on a gray background.

\mathcal{E}_{Text}

\mathcal{E} : Encoder

```

1 class CLIPTextEncoder(nn.Module):
2     def __init__(self):
3         super().__init__()
4         CLIP_id = "openai/clip-vit-base-patch32"
5         self.tokenizer = CLIPTokenizer.from_pretrained(CLIP_id)
6         self.text_encoder = CLIPTextModel.from_pretrained(CLIP_id)
7         self.device = "cuda" if torch.cuda.is_available() else "cpu"
8
9     for param in self.text_encoder.parameters():
10         param.requires_grad = False
11
12     self.text_encoder.eval()
13     self.text_encoder.to(self.device)
14
15     def forward(self, prompts):
16         inputs = self.tokenizer(
17             prompts,
18             padding="max_length",
19             truncation=True,
20             max_length=self.text_encoder.config.max_position_embeddings,
21             return_tensors="pt"
22         )
23         input_ids = inputs.input_ids.to(self.device)
24         attention_mask = inputs.attention_mask.to(self.device)
25
26         with torch.no_grad():
27             text_encoder_output = self.text_encoder(
28                 input_ids=input_ids,
29                 attention_mask=attention_mask
30             )
31         last_hidden_states = text_encoder_output.last_hidden_state
32
33         return last_hidden_states

```

Training SD from scratch

❖ DDMP

 β_t

```
1 class DDPMscheduler:
2     def __init__(
3         self,
4         random_generator,
5         train_timesteps=1000,
6         diffusion_beta_start=0.00085,
7         diffusion_beta_end=0.012
8     ):
9
10    self.betas = torch.linspace(
11        diffusion_beta_start ** 0.5, diffusion_beta_end ** 0.5,
12        train_timesteps, dtype=torch.float32) ** 2
13    self.alphas = 1.0 - self.betas
14    self.alphas_cumulative_product = torch.cumprod(self.alphas, dim=0)
15    self.one_val = torch.tensor(1.0)
16    self.prng_generator = random_generator
17    self.total_train_timesteps = train_timesteps
18    self.schedule_timesteps = torch.from_numpy(
19        np.arange(0, train_timesteps)[::-1].copy())
20
21    def set_steps(self, num_sampling_steps=50):
22        self.num_sampling_steps = num_sampling_steps
23        step_scaling_factor = self.total_train_timesteps // self.num_sampling_steps
24        timesteps_for_sampling = (
25            np.arange(0, num_sampling_steps) * step_scaling_factor
26            .round()[::-1].copy().astype(np.int64)
27        )
28        self.schedule_timesteps = torch.from_numpy(timesteps_for_sampling)
29
30    def _get_prior_timestep(self, current_timestep):
31        previous_t = current_timestep - self.total_train_timesteps
32        previous_t = previous_t // self.num_sampling_steps
33        return previous_t
```



Training SD from scratch

❖ DDMP

$$\alpha_t = 1 - \beta_t$$

```
1 class DDPMscheduler:
2     def __init__(
3         self,
4         random_generator,
5         train_timesteps=1000,
6         diffusion_beta_start=0.00085,
7         diffusion_beta_end=0.012
8     ):
9
10    self.betas = torch.linspace(
11        diffusion_beta_start ** 0.5, diffusion_beta_end ** 0.5,
12        train_timesteps, dtype=torch.float32) ** 2
13    self.alphas = 1.0 - self.betas
14    self.alphas_cumulative_product = torch.cumprod(self.alphas, dim=0)
15    self.one_val = torch.tensor(1.0)
16    self.prng_generator = random_generator
17    self.total_train_timesteps = train_timesteps
18    self.schedule_timesteps = torch.from_numpy(
19        np.arange(0, train_timesteps)[::-1].copy())
20
21    def set_steps(self, num_sampling_steps=50):
22        self.num_sampling_steps = num_sampling_steps
23        step_scaling_factor = self.total_train_timesteps // self.num_sampling_steps
24        timesteps_for_sampling = (
25            np.arange(0, num_sampling_steps) * step_scaling_factor
26            .round()[::-1].copy().astype(np.int64)
27        )
28        self.schedule_timesteps = torch.from_numpy(timesteps_for_sampling)
29
30    def _get_prior_timestep(self, current_timestep):
31        previous_t = current_timestep - self.total_train_timesteps
32        previous_t = previous_t // self.num_sampling_steps
33        return previous_t
```



Training SD from scratch

❖ DDMP

$$\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$$

```
1 class DDPMscheduler:
2     def __init__(
3         self,
4         random_generator,
5         train_timesteps=1000,
6         diffusion_beta_start=0.00085,
7         diffusion_beta_end=0.012
8     ):
9
10    self.betas = torch.linspace(
11        diffusion_beta_start ** 0.5, diffusion_beta_end ** 0.5,
12        train_timesteps, dtype=torch.float32) ** 2
13    self.alphas = 1.0 - self.betas
14    self.alphas_cumulative_product = torch.cumprod(self.alphas, dim=0)
15    self.one_val = torch.tensor(1.0)
16    self.prng_generator = random_generator
17    self.total_train_timesteps = train_timesteps
18    self.schedule_timesteps = torch.from_numpy(
19        np.arange(0, train_timesteps)[::-1].copy())
20
21    def set_steps(self, num_sampling_steps=50):
22        self.num_sampling_steps = num_sampling_steps
23        step_scaling_factor = self.total_train_timesteps // self.num_sampling_steps
24        timesteps_for_sampling = (
25            np.arange(0, num_sampling_steps) * step_scaling_factor
26            .round()[::-1].copy().astype(np.int64)
27        )
28        self.schedule_timesteps = torch.from_numpy(timesteps_for_sampling)
29
30    def _get_prior_timestep(self, current_timestep):
31        previous_t = current_timestep - self.total_train_timesteps
32        previous_t = previous_t // self.num_sampling_steps
33        return previous_t
```

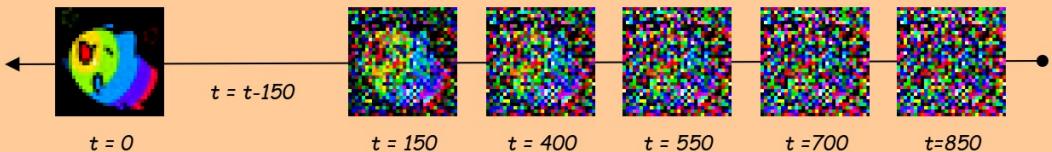
Training SD from scratch

❖ DDMP

Reverse Diffusion (Denoising)



Reverse Diffusion (Denoising)



```
1 class DDPMscheduler:
2     def __init__(self,
3                  random_generator,
4                  train_timesteps=1000,
5                  diffusion_beta_start=0.00085,
6                  diffusion_beta_end=0.012
7                  ):
8
9
10    self.betas = torch.linspace(
11        diffusion_beta_start ** 0.5, diffusion_beta_end ** 0.5,
12        train_timesteps, dtype=torch.float32) ** 2
13    self.alphas = 1.0 - self.betas
14    self.alphas_cumulative_product = torch.cumprod(self.alphas, dim=0)
15    self.one_val = torch.tensor(1.0)
16    self.prng_generator = random_generator
17    self.total_train_timesteps = train_timesteps
18    self.schedule_timesteps = torch.from_numpy(
19        np.arange(0, train_timesteps)[::-1].copy())
20
21    def set_steps(self, num_sampling_steps=50):
22        self.num_sampling_steps = num_sampling_steps
23        step_scaling_factor = self.total_train_timesteps // self.num_sampling_steps
24        timesteps_for_sampling = (
25            np.arange(0, num_sampling_steps) * step_scaling_factor
26            .round()[::-1].copy().astype(np.int64))
27        self.schedule_timesteps = torch.from_numpy(timesteps_for_sampling)
28
29    def _get_prior_timestep(self, current_timestep):
30        previous_t = current_timestep - self.total_train_timesteps
31        previous_t = previous_t // self.num_sampling_steps
32        return previous_t
```

Training SD from scratch

❖ DDMP

$$\sigma_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

```
34 def _calculate_variance(self, timestep):
35     prev_t = self._get_prior_timestep(timestep)
36     alpha_cumprod_t = self.alphas_cumulative_product[timestep]
37     alpha_cumprod_t_prev = self.alphas_cumulative_product[prev_t] if prev_t >= 0 else self.one_val
38     beta_t_current = 1 - alpha_cumprod_t / alpha_cumprod_t_prev
39     variance_value = (1 - alpha_cumprod_t_prev) / (1 - alpha_cumprod_t) * beta_t_current
40     variance_value = torch.clamp(variance_value, min=1e-20)
41     return variance_value
```

Training SD from scratch

❖ DDMP

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta)$$

```
43 def step(self, current_t, current_latents, model_prediction):
44     t = current_t
45     prev_t = self._get_prior_timestep(t)
46
47     alpha_cumprod_t = self.alphas_cumulative_product[t]
48     alpha_cumprod_t_prev = self.alphas_cumulative_product[prev_t] if prev_t >= 0 else self.one_val
49     beta_cumprod_t = 1 - alpha_cumprod_t
50     beta_cumprod_t_prev = 1 - alpha_cumprod_t_prev
51     alpha_t_current = alpha_cumprod_t / alpha_cumprod_t_prev
52     beta_t_current = 1 - alpha_t_current
53
54     predicted_ori = (
55         cur_latents - beta_cumprod_t ** 0.5 * model_prediction
56     ) / alpha_cumprod_t ** 0.5
57
58     ori_coeff = (alpha_cumprod_t_prev ** 0.5 * beta_t_current) / beta_cumprod_t
59     cur_coeff = alpha_t_current ** 0.5 * beta_cumprod_t_prev / beta_cumprod_t
60
61     predicted_prior_mean = ori_coeff * predicted_ori + cur_coeff * cur_latents
62
63     variance_term = 0
64     if t > 0:
65         target_device = model_prediction.device
66         noise_component = torch.randn(
67             model_prediction.shape,
68             generator=self.prng_generator,
69             device=target_device,
70             dtype=model_prediction.dtype
71         )
72         variance_term = (self._calculate_variance(t) ** 0.5) * noise_component
73
74     predicted_prior_sample = predicted_prior_mean + variance_term
75     return predicted_prior_sample
```

Training SD from scratch

❖ DDMP

$$\mu_t = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} x_0$$

$$\alpha_t = 1 - \beta_t$$

$$\beta_t = 1 - \alpha_t$$

```

43 def step(self, current_t, current_latents, model_prediction):
44     t = current_t
45     prev_t = self._get_prior_timestep(t)
46
47     alpha_cumprod_t = self.alphas_cumulative_product[t]
48     alpha_cumprod_t_prev = self.alphas_cumulative_product[prev_t] if prev_t >= 0 else self.one_val
49     beta_cumprod_t = 1 - alpha_cumprod_t
50     beta_cumprod_t_prev = 1 - alpha_cumprod_t_prev
51     alpha_t_current = alpha_cumprod_t / alpha_cumprod_t_prev
52     beta_t_current = 1 - alpha_t_current
53
54     predicted_ori = (
55         cur_latents - beta_cumprod_t ** 0.5 * model_prediction
56     ) / alpha_cumprod_t ** 0.5
57
58     ori_coeff = (alpha_cumprod_t_prev ** 0.5 * beta_t_current) / beta_cumprod_t
59     cur_coeff = alpha_t_current ** 0.5 * beta_cumprod_t_prev / beta_cumprod_t
60
61     predicted_prior_mean = ori_coeff * predicted_ori + cur_coeff * cur_latents
62
63     variance_term = 0
64     if t > 0:
65         target_device = model_prediction.device
66         noise_component = torch.randn(
67             model_prediction.shape,
68             generator=self.prng_generator,
69             device=target_device,
70             dtype=model_prediction.dtype
71         )
72         variance_term = (self._calculate_variance(t) ** 0.5) * noise_component
73
74     predicted_prior_sample = predicted_prior_mean + variance_term
75     return predicted_prior_sample

```

Training SD from scratch

❖ DDMP

$$\sqrt{\beta_t} \epsilon$$

```
43 def step(self, current_t, current_latents, model_prediction):
44     t = current_t
45     prev_t = self._get_prior_timestep(t)
46
47     alpha_cumprod_t = self.alphas_cumulative_product[t]
48     alpha_cumprod_t_prev = self.alphas_cumulative_product[prev_t] if prev_t >= 0 else self.one_val
49     beta_cumprod_t = 1 - alpha_cumprod_t
50     beta_cumprod_t_prev = 1 - alpha_cumprod_t_prev
51     alpha_t_current = alpha_cumprod_t / alpha_cumprod_t_prev
52     beta_t_current = 1 - alpha_t_current
53
54     predicted_ori = (
55         cur_latents - beta_cumprod_t ** 0.5 * model_prediction
56     ) / alpha_cumprod_t ** 0.5
57
58     ori_coeff = (alpha_cumprod_t_prev ** 0.5 * beta_t_current) / beta_cumprod_t
59     cur_coeff = alpha_t_current ** 0.5 * beta_cumprod_t_prev / beta_cumprod_t
60
61     predicted_prior_mean = ori_coeff * predicted_ori + cur_coeff * cur_latents
62
63     variance_term = 0
64     if t > 0:
65         target_device = model_prediction.device
66         noise_component = torch.randn(
67             model_prediction.shape,
68             generator=self.prng_generator,
69             device=target_device,
70             dtype=model_prediction.dtype
71         )
72         variance_term = (self._calculate_variance(t) ** 0.5) * noise_component
73
74     predicted_prior_sample = predicted_prior_mean + variance_term
75     return predicted_prior_sample
```

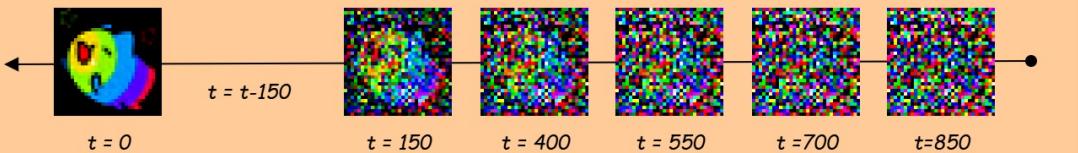
Training SD from scratch

❖ DDMP

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta \right) + \sqrt{\beta_t} \epsilon$$

$$\epsilon \sim \mathcal{N}(0,1)$$

Reverse Diffusion (Denoising)



```

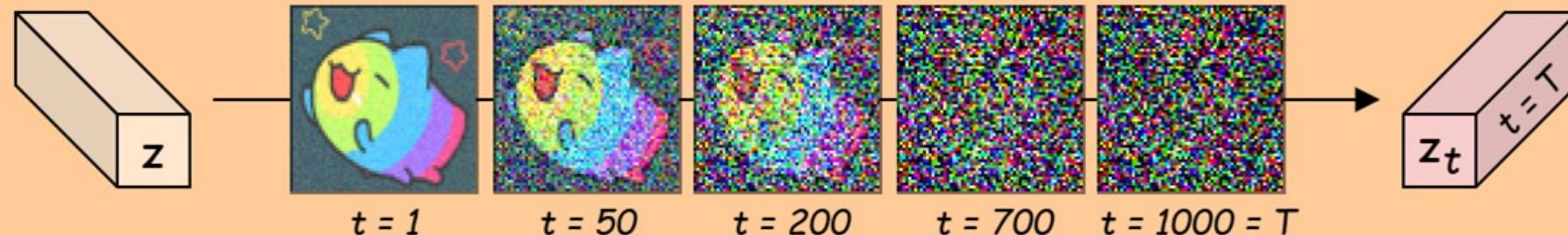
43 def step(self, current_t, current_latents, model_prediction):
44     t = current_t
45     prev_t = self._get_prior_timestep(t)
46
47     alpha_cumprod_t = self.alphas_cumulative_product[t]
48     alpha_cumprod_t_prev = self.alphas_cumulative_product[prev_t] if prev_t >= 0 else self.one_val
49     beta_cumprod_t = 1 - alpha_cumprod_t
50     beta_cumprod_t_prev = 1 - alpha_cumprod_t_prev
51     alpha_t_current = alpha_cumprod_t / alpha_cumprod_t_prev
52     beta_t_current = 1 - alpha_t_current
53
54     predicted_ori = (
55         cur_latents - beta_cumprod_t ** 0.5 * model_prediction
56     ) / alpha_cumprod_t ** 0.5
57
58     ori_coeff = (alpha_cumprod_t_prev ** 0.5 * beta_t_current) / beta_cumprod_t
59     cur_coeff = alpha_t_current ** 0.5 * beta_cumprod_t_prev / beta_cumprod_t
60
61     predicted_prior_mean = ori_coeff * predicted_ori + cur_coeff * cur_latents
62
63     variance_term = 0
64     if t > 0:
65         target_device = model_prediction.device
66         noise_component = torch.randn(
67             model_prediction.shape,
68             generator=self.prng_generator,
69             device=target_device,
70             dtype=model_prediction.dtype
71         )
72         variance_term = (self._calculate_variance(t) ** 0.5) * noise_component
73
74         predicted_prior_sample = predicted_prior_mean + variance_term
75     return predicted_prior_sample
    
```

Training SD from scratch

❖ DDMP

Latent Space

Forward Diffusion (The images below are used to illustrate adding noise to the latent z)



$$q(x_t \mid x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I) = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t$$

Training SD from scratch

❖ DDMP

$$q(x_t \mid x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I) = \boxed{\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t}$$

```
77 def add_noise(self, initial_samples, noise_timesteps):
78     alphas_cumprod = self.alphas_cumulative_product.to(
79         device=initial_samples.device,
80         dtype=initial_samples.dtype
81     )
82     noise_timesteps = noise_timesteps.to(initial_samples.device)
83     sqrt_alpha_cumprod = alphas_cumprod=noise_timesteps] ** 0.5
84     sqrt_alpha_cumprod = sqrt_alpha_cumprod.view(
85         sqrt_alpha_cumprod.shape[0], *[1] * (initial_samples.ndim - 1))
86     )
87     sqrt_one_minus_alpha_cumprod = (1 - alphas_cumprod=noise_timesteps]) ** 0.5
88     sqrt_one_minus_alpha_cumprod = sqrt_one_minus_alpha_cumprod.view(
89         sqrt_one_minus_alpha_cumprod.shape[0], *[1] * (initial_samples.ndim - 1))
90     )
91     random_noise = torch.randn(
92         initial_samples.shape, generator=self.prng_generator,
93         device=initial_samples.device, dtype=initial_samples.dtype
94     )
95     a = sqrt_alpha_cumprod * initial_samples
96     b = sqrt_one_minus_alpha_cumprod * random_noise
97     noisy_result = a + b
98     return noisy_result, random_noise
```

Training SD from scratch

❖ DDMP

$$q(x_t \mid x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I) = \sqrt{\bar{\alpha}_t}x_0 + \boxed{\sqrt{1 - \bar{\alpha}_t}\epsilon_t}$$

```
77 def add_noise(self, initial_samples, noise_timesteps):
78     alphas_cumprod = self.alphas_cumulative_product.to(
79         device=initial_samples.device,
80         dtype=initial_samples.dtype
81     )
82     noise_timesteps = noise_timesteps.to(initial_samples.device)
83     sqrt_alpha_cumprod = alphas_cumprod=noise_timesteps] ** 0.5
84     sqrt_alpha_cumprod = sqrt_alpha_cumprod.view(
85         sqrt_alpha_cumprod.shape[0], *[1] * (initial_samples.ndim - 1))
86     )
87     sqrt_one_minus_alpha_cumprod = (1 - alphas_cumprod=noise_timesteps]) ** 0.5
88     sqrt_one_minus_alpha_cumprod = sqrt_one_minus_alpha_cumprod.view(
89         sqrt_one_minus_alpha_cumprod.shape[0], *[1] * (initial_samples.ndim - 1))
90     )
91     random_noise = torch.randn(
92         initial_samples.shape, generator=self.prng_generator,
93         device=initial_samples.device, dtype=initial_samples.dtype
94     )
95     a = sqrt_alpha_cumprod * initial_samples
96     b = sqrt_one_minus_alpha_cumprod * random_noise
97     noisy_result = a + b
98     return noisy_result, random_noise
```

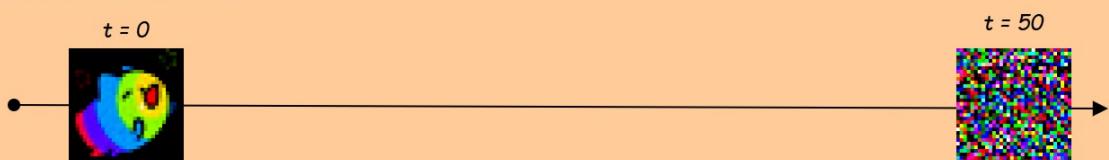
Training SD from scratch

❖ DDMP

Forward Diffusion



Forward Diffusion

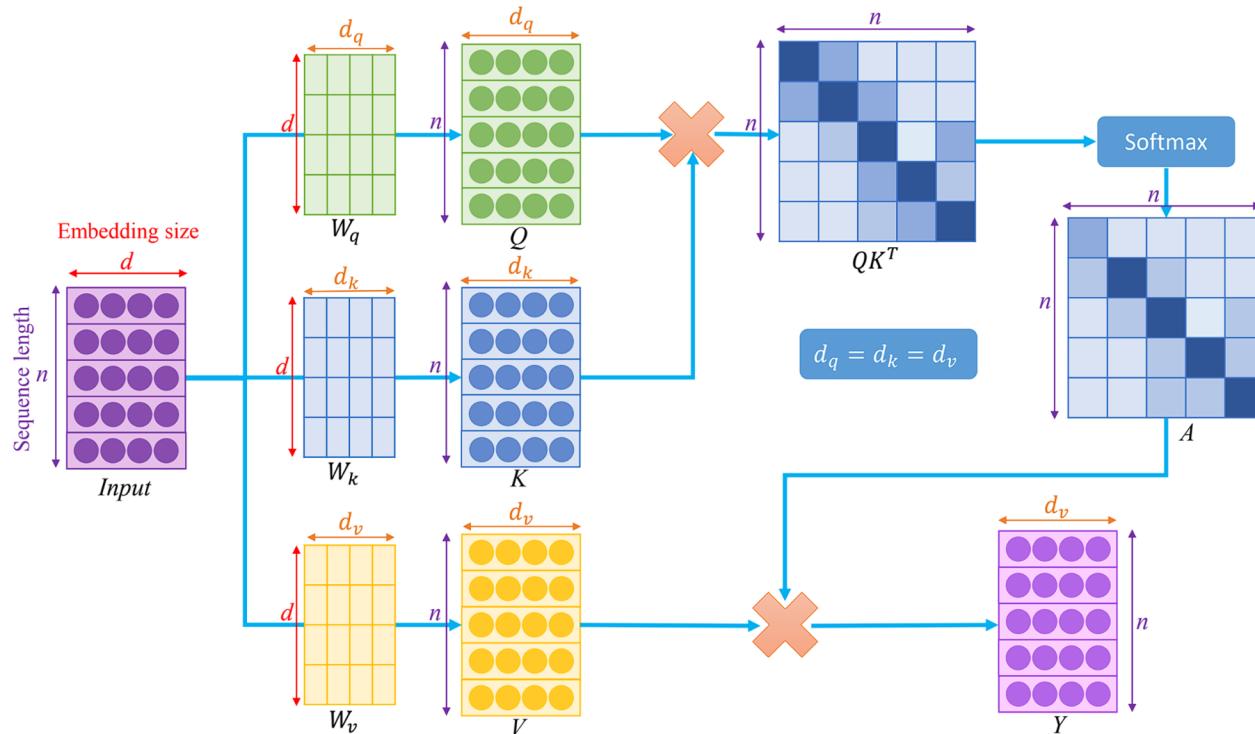


$$\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t$$

```
77 def add_noise(self, initial_samples, noise_timesteps):
78     alphas_cumprod = self.alphas_cumulative_product.to(
79         device=initial_samples.device,
80         dtype=initial_samples.dtype
81     )
82     noise_timesteps = noise_timesteps.to(initial_samples.device)
83     sqrt_alpha_cumprod = alphas_cumprod=noise_timesteps ** 0.5
84     sqrt_alpha_cumprod = sqrt_alpha_cumprod.view(
85         sqrt_alpha_cumprod.shape[0], *[1] * (initial_samples.ndim - 1))
86     sqrt_one_minus_alpha_cumprod = (1 - alphas_cumprod=noise_timesteps) ** 0.5
87     sqrt_one_minus_alpha_cumprod = sqrt_one_minus_alpha_cumprod.view(
88         sqrt_one_minus_alpha_cumprod.shape[0], *[1] * (initial_samples.ndim - 1))
89
90     random_noise = torch.randn(
91         initial_samples.shape, generator=self.prng_generator,
92         device=initial_samples.device, dtype=initial_samples.dtype
93     )
94     a = sqrt_alpha_cumprod * initial_samples
95     b = sqrt_one_minus_alpha_cumprod * random_noise
96     noisy_result = a + b
97     return noisy_result, random_noise
98
```

Training SD from scratch

❖ Self-Attention Review



$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

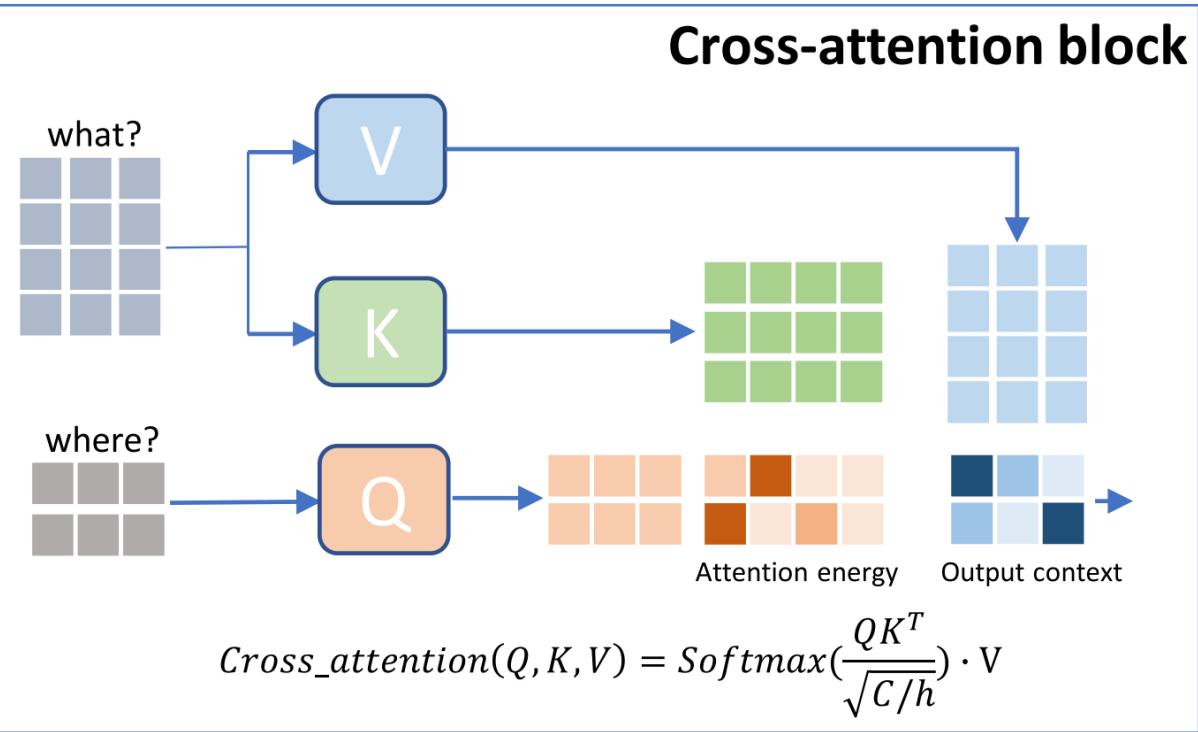
```

1 class SelfAttention(nn.Module):
2     def __init__(self, num_attn_heads, hidden_dim,
3                  in_proj_bias=True, out_proj_bias=True):
4         super().__init__()
5         self.num_heads = num_attn_heads
6         self.head_size = hidden_dim // num_attn_heads
7
8         self.qkv_proj = nn.Linear(hidden_dim, 3 * hidden_dim, bias=in_proj_bias)
9         self.output_proj = nn.Linear(hidden_dim, hidden_dim, bias=out_proj_bias)
10
11    def forward(self, features, use_causal_mask=False):
12        b, s, d = features.shape
13
14        qkv_combined = self.qkv_proj(features)
15        q_mat, k_mat, v_mat = torch.chunk(qkv_combined, 3, dim=-1)
16
17        q_mat = q_mat.view(b, s, self.num_heads, self.head_size).permute(0, 2, 1, 3)
18        k_mat = k_mat.view(b, s, self.num_heads, self.head_size).permute(0, 2, 1, 3)
19        v_mat = v_mat.view(b, s, self.num_heads, self.head_size).permute(0, 2, 1, 3)
20
21        qk = torch.matmul(q_mat, k_mat.transpose(-2, -1))
22        sqrt_qk = qk / math.sqrt(self.head_size)
23
24        if use_causal_mask:
25            causal_mask = torch.triu(
26                torch.ones_like(sqrt_qk, dtype=torch.bool), diagonal=1)
27            sqrt_qk = sqrt_qk.masked_fill(causal_mask, -torch.inf)
28
29        attn_weights = torch.softmax(sqrt_qk, dim=-1)
30        attn_values = torch.matmul(attn_weights, v_mat)
31
32        attn_values = attn_values.permute(0, 2, 1, 3).contiguous()
33        attn_values = attn_values.view(b, s, d)
34
35        final_output = self.output_proj(attn_values)
36        return final_output

```

Training SD from scratch

❖ Cross-Attention Review



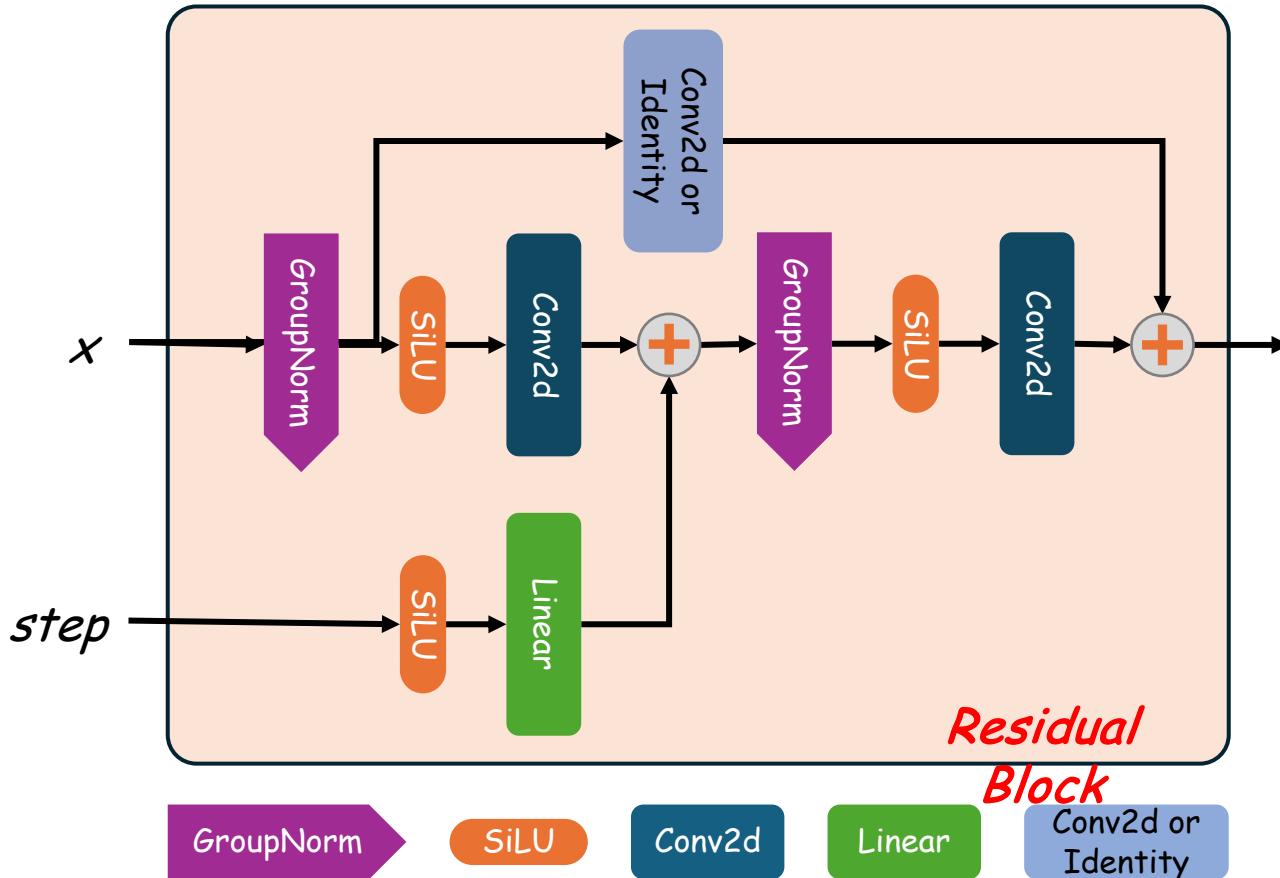
```

1 class CrossAttention(nn.Module):
2     def __init__(self, num_attn_heads, query_dim, context_dim,
3                  in_proj_bias=True, out_proj_bias=True):
4         super().__init__()
5         self.num_heads = num_attn_heads
6         self.head_size = query_dim // num_attn_heads
7
8         self.query_map = nn.Linear(query_dim, query_dim, bias=in_proj_bias)
9         self.key_map = nn.Linear(context_dim, query_dim, bias=in_proj_bias)
10        self.value_map = nn.Linear(context_dim, query_dim, bias=in_proj_bias)
11
12        self.output_map = nn.Linear(query_dim, query_dim, bias=out_proj_bias)
13
14    def forward(self, query_input, context_input):
15        b_q, s_q, d_q = query_input.shape
16        _, s_kv, _ = context_input.shape
17
18        q_mat = self.query_map(query_input)
19        k_mat = self.key_map(context_input)
20        v_mat = self.value_map(context_input)
21
22        q_mat = q_mat.view(b_q, s_q, self.num_heads, self.head_size).permute(0, 2, 1, 3)
23        k_mat = k_mat.view(b_q, s_kv, self.num_heads, self.head_size).permute(0, 2, 1, 3)
24        v_mat = v_mat.view(b_q, s_kv, self.num_heads, self.head_size).permute(0, 2, 1, 3)
25
26        qk = torch.matmul(q_mat, k_mat.transpose(-2, -1))
27        sqrt_qk = qk / math.sqrt(self.head_size)
28        attn_weights = torch.softmax(sqrt_qk, dim=-1)
29
30        attn_values = torch.matmul(attn_weights, v_mat)
31        attn_values = attn_values.permute(0, 2, 1, 3).contiguous()
32        attn_values = attn_values.view(b_q, s_q, d_q)
33
34        final_output = self.output_map(attn_values)
35        return final_output

```

Training SD from scratch

❖ Unet – Residual Block



```

1 class UNET_ResidualBlock(nn.Module):
2     def __init__(self, in_channels, out_channels, time_dim=1280):
3         super().__init__()
4         self.gn_feature = nn.GroupNorm(32, in_channels)
5         self.conv_feature = nn.Conv2d(
6             in_channels, out_channels, kernel_size=3, padding=1)
7         self.time_embedding_proj = nn.Linear(time_dim, out_channels)
8
9         self.gn_merged = nn.GroupNorm(32, out_channels)
10        self.conv_merged = nn.Conv2d(
11            out_channels, out_channels, kernel_size=3, padding=1)
12
13        if in_channels == out_channels:
14            self.residual_connection = nn.Identity()
15        else:
16            self.residual_connection = nn.Conv2d(
17                in_channels, out_channels, kernel_size=1, padding=0)
18
19    def forward(self, input_feature, time_emb):
20        residual = input_feature
21
22        h = self.gn_feature(input_feature)
23        h = F.silu(h)
24        h = self.conv_feature(h)
25
26        time_emb_processed = F.silu(time_emb)
27        time_emb_projected = self.time_embedding_proj(time_emb_processed)
28        time_emb_projected = time_emb_projected.unsqueeze(-1).unsqueeze(-1)
29
30        merged_feature = h + time_emb_projected
31        merged_feature = self.gn_merged(merged_feature)
32        merged_feature = F.silu(merged_feature)
33        merged_feature = self.conv_merged(merged_feature)
34
35        output = merged_feature + self.residual_connection(residual)
36        return output

```

Training SD from scratch

❖ Unet – Attention Block

```

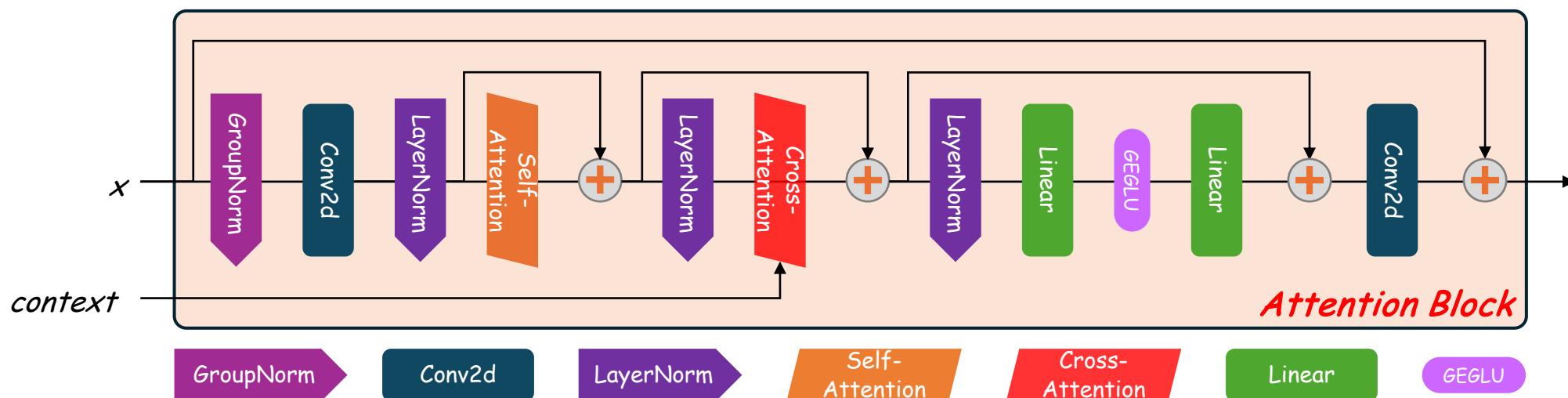
1 class UNET_AttentionBlock(nn.Module):
2     def __init__(self, num_heads, head_dim, context_dim=512):
3         super().__init__()
4         embed_dim = num_heads * head_dim
5
6         self.gn_in = nn.GroupNorm(32, embed_dim, eps=1e-6)
7         self.proj_in = nn.Conv2d(embed_dim, embed_dim, kernel_size=1, padding=0)
8
9         self.ln_1 = nn.LayerNorm(embed_dim)
10        self.attn_1 = SelfAttention(num_heads, embed_dim, in_proj_bias=False)
11        self.ln_2 = nn.LayerNorm(embed_dim)
12        self.attn_2 = CrossAttention(
13            num_heads, embed_dim, context_dim, in_proj_bias=False)
14        self.ln_3 = nn.LayerNorm(embed_dim)
15
16        self.ffn_gelu = nn.Linear(embed_dim, 4 * embed_dim * 2)
17        self.ffn_out = nn.Linear(4 * embed_dim, embed_dim)
18        self.proj_out = nn.Conv2d(embed_dim, embed_dim, kernel_size=1, padding=0)

```

```

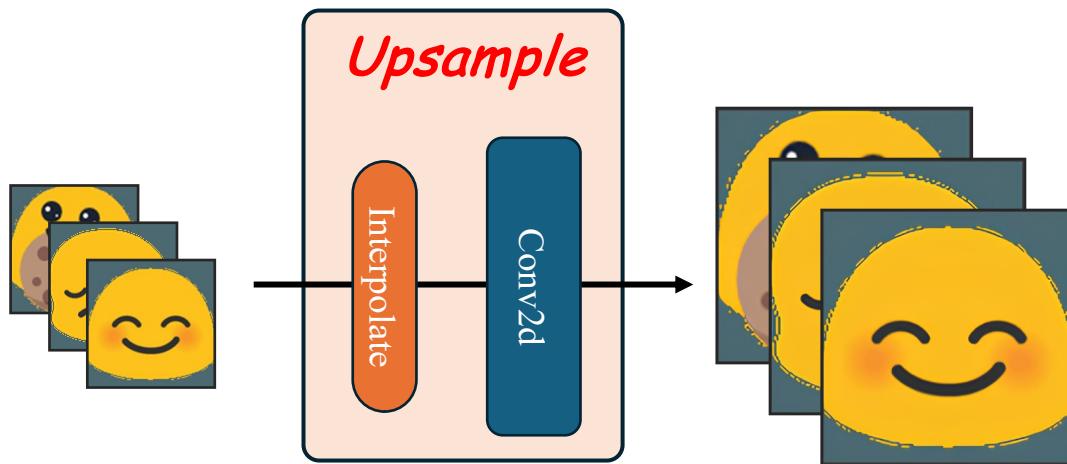
20     def forward(self, input_tensor, context_tensor):
21         skip_connection = input_tensor
22         B, C, H, W = input_tensor.shape
23         HW = H * W
24
25         h = self.proj_in(self.gn_in(input_tensor))
26         h = h.view(B, C, HW).transpose(-1, -2)
27
28         attn1_skip = h
29         h = self.attn_1(self.ln_1(h)) + attn1_skip
30
31         attn2_skip = h
32         h = self.attn_2(self.ln_2(h), context_tensor) + attn2_skip
33
34         ffn_skip = h
35         intermediate, gate = self.ffn_gelu(self.ln_3(h)).chunk(2, dim=-1)
36         h = self.ffn_out(intermediate * F.gelu(gate)) + ffn_skip
37
38         h = h.transpose(-1, -2).view(B, C, H, W)
39         output = self.proj_out(h) + skip_connection
40         return output

```



Training SD from scratch

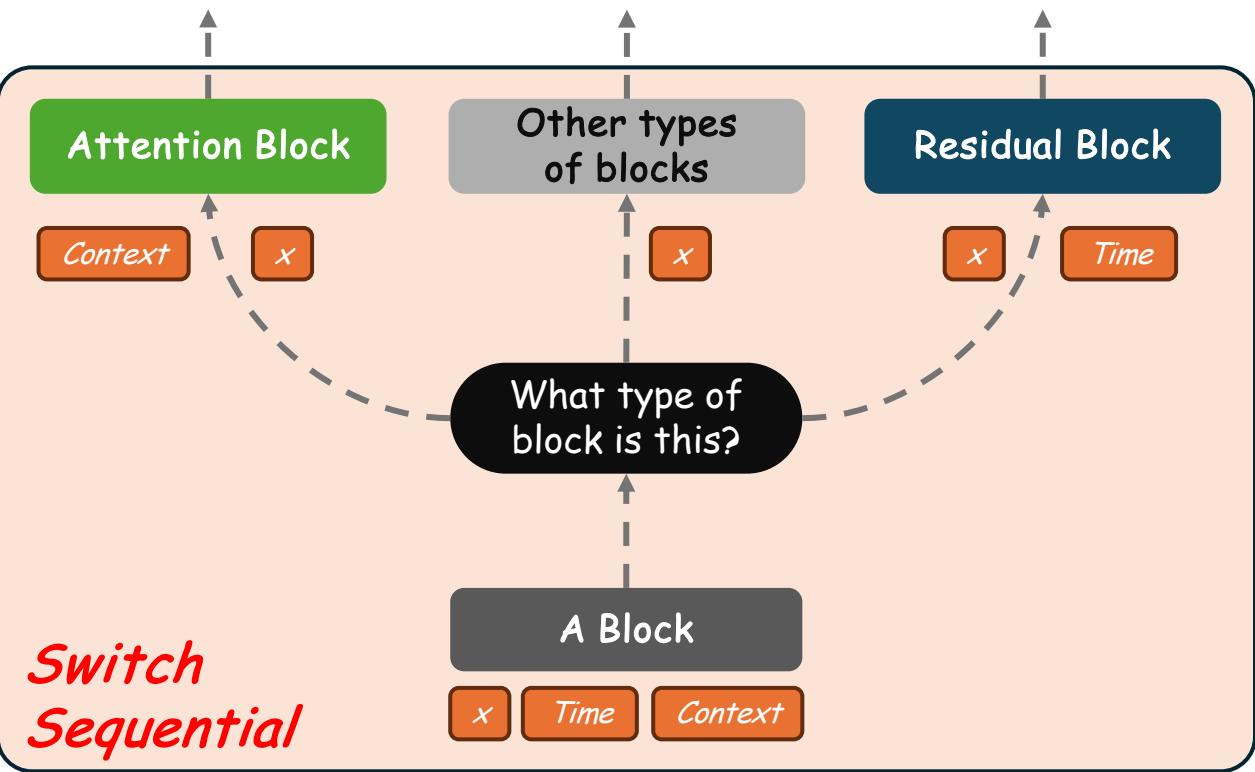
❖ Unet – Upsample Block



```
1 class Upsample(nn.Module):
2     def __init__(self, num_channels):
3         super().__init__()
4         self.conv = nn.Conv2d(
5             num_channels, num_channels,
6             kernel_size=3, padding=1)
7
8     def forward(self, feature_map):
9         x = F.interpolate(feature_map,
10                          scale_factor=2,
11                          mode='nearest')
12         x = self.conv(x)
13
14     return x
```

Training SD from scratch

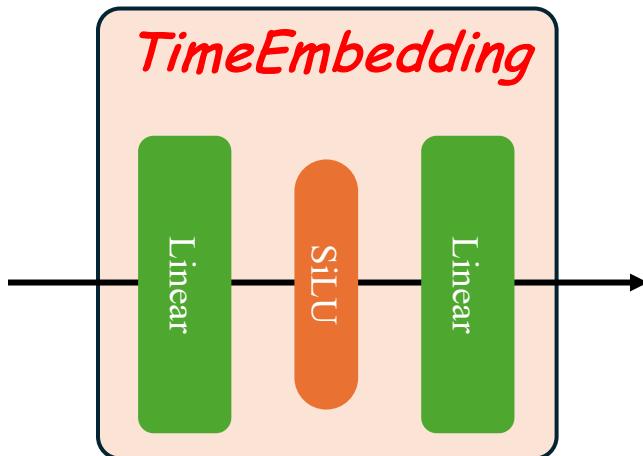
❖ Unet – Switch Sequential Block



```
16 class SwitchSequential(nn.Sequential):
17     def forward(self, x, guidance_context, time_embedding):
18         for module_instance in self:
19             if isinstance(module_instance, UNET_AttentionBlock):
20                 x = module_instance(x, guidance_context)
21             elif isinstance(module_instance, UNET_ResidualBlock):
22                 x = module_instance(x, time_embedding)
23             else:
24                 x = module_instance(x)
25         return x
```

Training SD from scratch

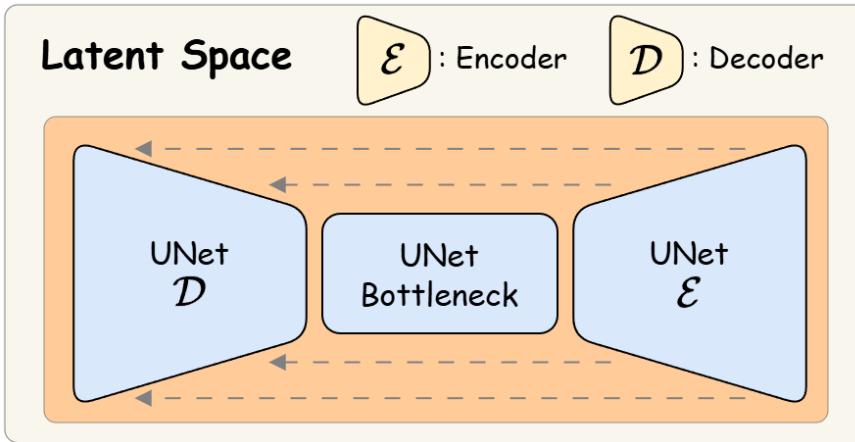
❖ Unet – TimeEmbedding Block



```
24 class TimeEmbedding(nn.Module):  
25     def __init__(self, n_embd):  
26         super().__init__()  
27         self.proj1 = nn.Linear(n_embd, 4 * n_embd)  
28         self.proj2 = nn.Linear(4 * n_embd, 4 * n_embd)  
29  
30     def forward(self, x):  
31         x = self.proj1(x)  
32         x = F.silu(x)  
33         x = self.proj2(x)  
34         return x
```

Training SD from scratch

❖ Unet



```

1 class UNET(nn.Module):
2     def __init__(self, h_dim, n_head):
3         super().__init__()
4
5         self.down_blocks = nn.ModuleList([
6             SwitchSequential(nn.Conv2d(4, h_dim,
7                 kernel_size=3, padding=1)),
8             SwitchSequential(
9                 UNET_ResidualBlock(h_dim, h_dim),
10                UNET_AttentionBlock(n_head, (h_dim)//n_head)),
11             SwitchSequential(
12                 UNET_ResidualBlock(h_dim, h_dim),
13                UNET_AttentionBlock(n_head, (h_dim)//n_head)),
14             SwitchSequential(nn.Conv2d(h_dim, h_dim,
15                 kernel_size=3, stride=2, padding=1)),
16             SwitchSequential(UNET_ResidualBlock(h_dim, 2*h_dim)),
17             SwitchSequential(UNET_ResidualBlock(2*h_dim, 2*h_dim)),
18         ])

```

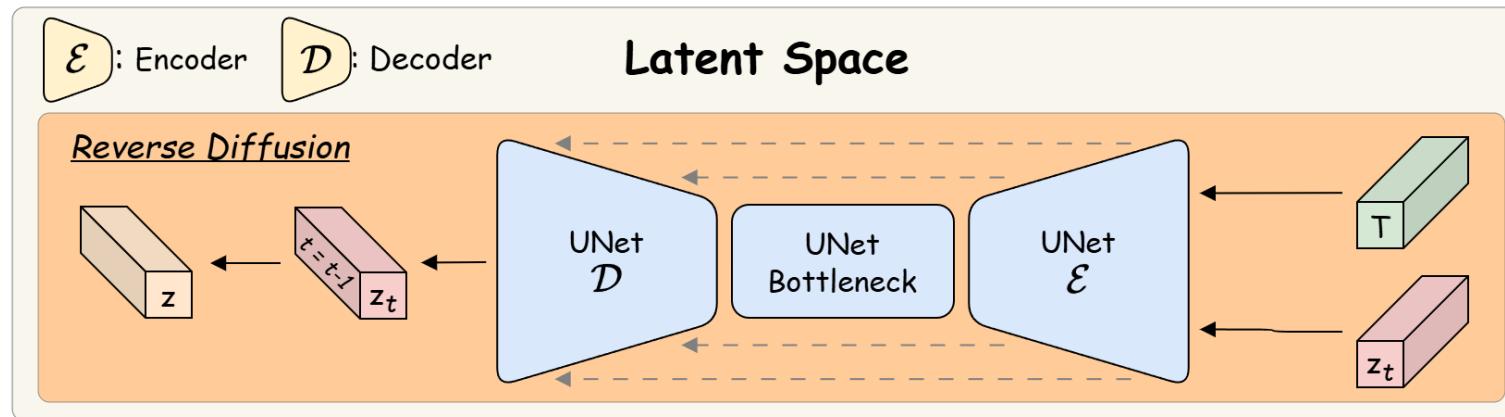
```

20     self.mid_block = SwitchSequential(
21         UNET_ResidualBlock(2*h_dim, 2*h_dim),
22         UNET_AttentionBlock(n_head, (2*h_dim)//n_head),
23         UNET_ResidualBlock(2*h_dim, 2*h_dim),
24     )
25
26     self.up_blocks = nn.ModuleList([
27         SwitchSequential(UNET_ResidualBlock(2*h_dim + 2*h_dim, 2*h_dim)),
28         SwitchSequential(UNET_ResidualBlock(2*h_dim + 2*h_dim, 2*h_dim)),
29         SwitchSequential(
30             UNET_ResidualBlock(2*h_dim + h_dim, 2*h_dim),
31             Upsample(2*h_dim)),
32             SwitchSequential(
33                 UNET_ResidualBlock(2*h_dim + h_dim, h_dim),
34                 UNET_AttentionBlock(n_head, h_dim//n_head)),
35             SwitchSequential(
36                 UNET_ResidualBlock(2*h_dim, h_dim),
37                 UNET_AttentionBlock(n_head, h_dim//n_head)),
38             SwitchSequential(
39                 UNET_ResidualBlock(2*h_dim, h_dim),
40                 UNET_AttentionBlock(n_head, h_dim//n_head)),
41         ])

```

Training SD from scratch

❖ Unet



```

43 def forward(self, latent_input, context_embedding, time_embedding):
44     down_block_residuals = []
45     current_feature_map = latent_input
46
47     for block in self.down_blocks:
48         current_feature_map = block(
49             current_feature_map, context_embedding, time_embedding)
50         down_block_residuals.append(current_feature_map)
51
52     current_feature_map = self.mid_block(
53         current_feature_map, context_embedding, time_embedding)
54
55     for block in self.up_blocks:
56         residual = down_block_residuals.pop()
57         current_feature_map = torch.cat(
58             (current_feature_map, residual), dim=1)
59         current_feature_map = block(
60             current_feature_map, context_embedding, time_embedding)
61
62     return current_feature_map
  
```

```

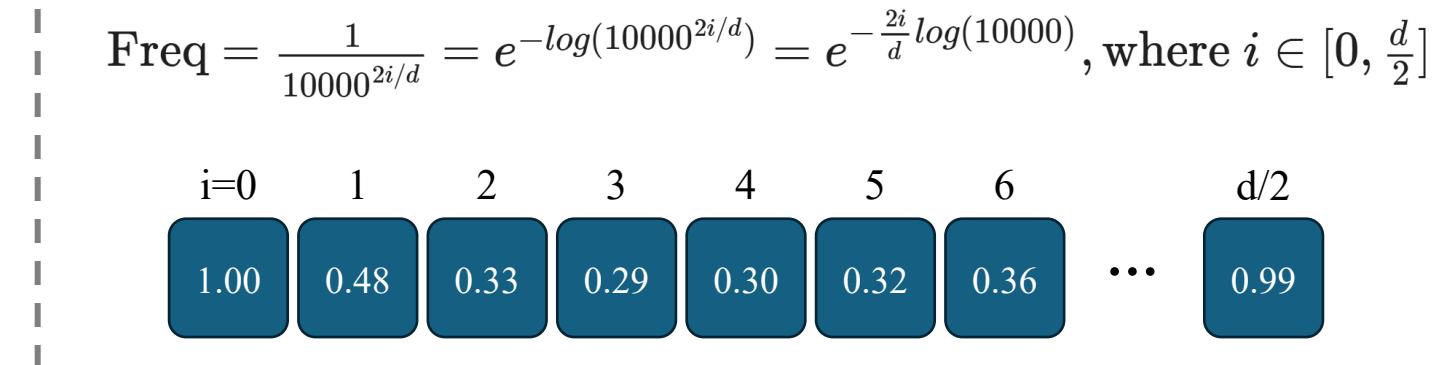
61 class UNETOutputLayer(nn.Module):
62     def __init__(self, input_channels, output_channels):
63         super().__init__()
64         self.final_groupnorm = nn.GroupNorm(32, input_channels)
65         self.final_conv = nn.Conv2d(input_channels, output_channels,
66                                   kernel_size=3, padding=1)
67
68     def forward(self, feature_map):
69         norm_map = self.final_groupnorm(feature_map)
70         activated_map = F.silu(norm_map)
71         output_map = self.final_conv(activated_map)
72
73     return output_map
  
```

Training SD from scratch

❖ Timestep Handling

fixed

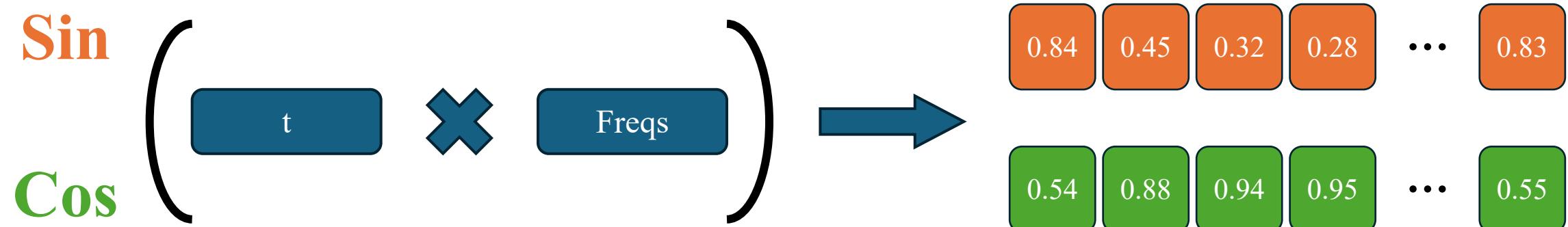
$$PE_{(pos,2i)} = \sin(pos / 10000^{2i/d_{\text{model}}})$$
$$PE_{(pos,2i+1)} = \cos(pos / 10000^{2i/d_{\text{model}}})$$



```
9 def embed_timesteps(timesteps, embedding_dim=320):  
10    half_dim = embedding_dim // 2  
11    freqs = torch.exp(-math.log(10000) *  
12                           torch.arange(half_dim, dtype=torch.float32) /  
13                           half_dim).to(device=timesteps.device)  
14    args = timesteps[:, None].float() * freqs[None, :]  
15    return torch.cat([torch.cos(args), torch.sin(args)], dim=-1)
```

Training SD from scratch

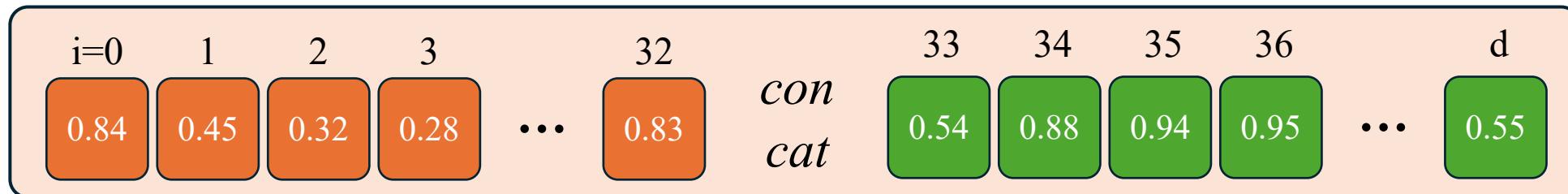
❖ Timestep Handling



```
9 def embed_timesteps(timesteps, embedding_dim=320):  
10    half_dim = embedding_dim // 2  
11    freqs = torch.exp(-math.log(10000) *  
12    torch.arange(half_dim, dtype=torch.float32) /  
13    half_dim).to(device=timesteps.device)  
14    args = timesteps[:, None].float() * freqs[None, :]  
15    return torch.cat([torch.cos(args), torch.sin(args)], dim=-1)
```

Training SD from scratch

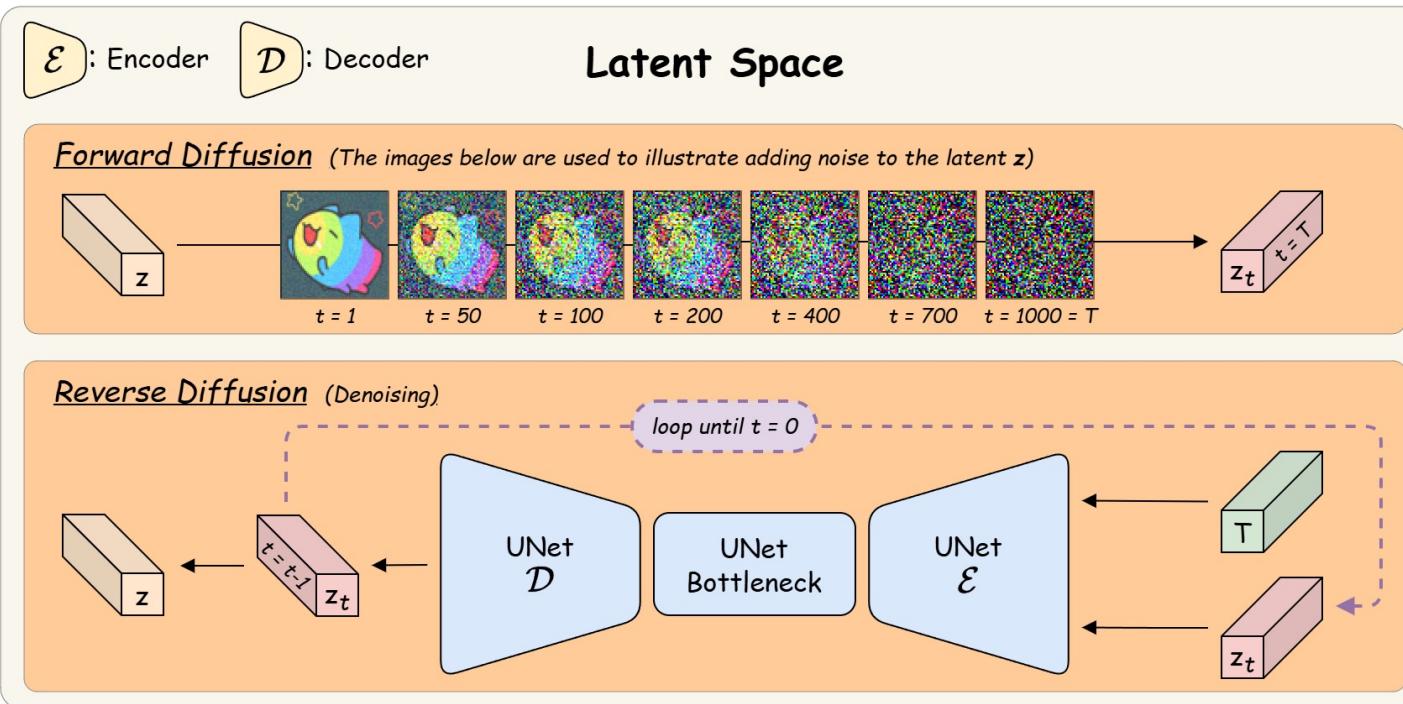
❖ Timestep Handling



```
9 def embed_timesteps(timesteps, embedding_dim=320):
10     half_dim = embedding_dim // 2
11     freqs = torch.exp(-math.log(10000) *
12                         torch.arange(half_dim, dtype=torch.float32) /
13                         half_dim).to(device=timesteps.device)
14     args = timesteps[:, None].float() * freqs[None, :]
15     return torch.cat([torch.cos(args), torch.sin(args)], dim=-1)
```

Training SD from scratch

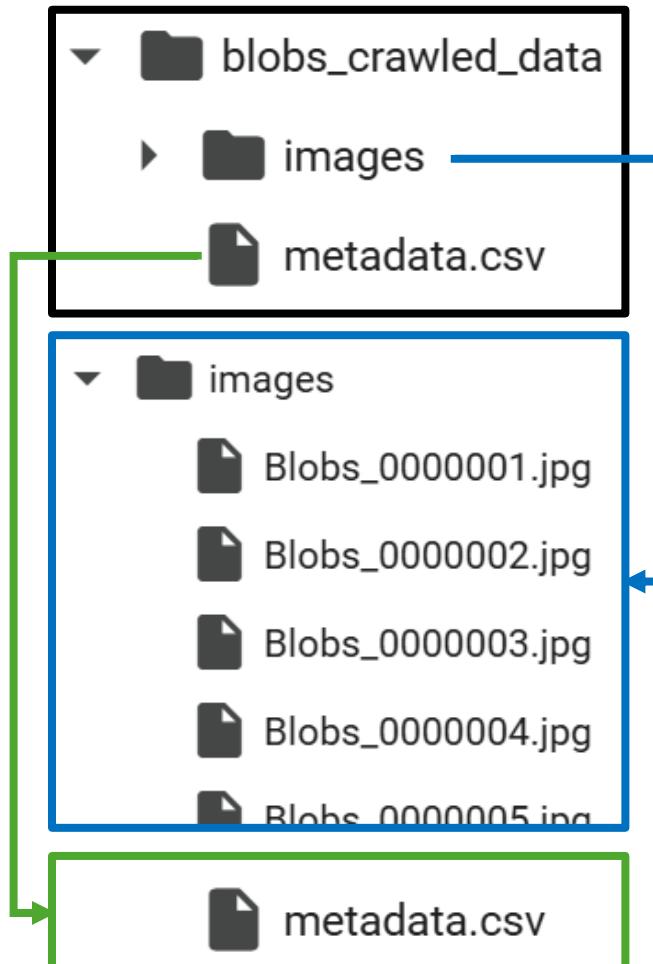
❖ Diffusion Process



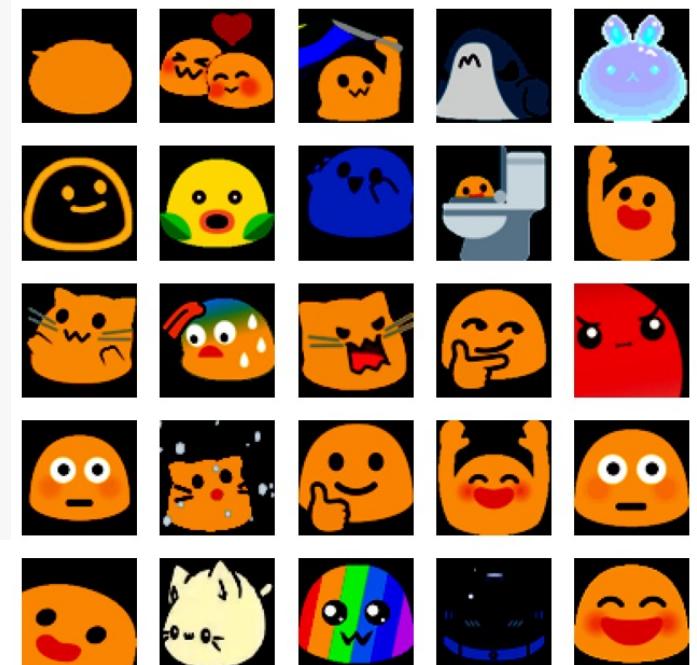
```
1 class Diffusion(nn.Module):
2     def __init__(self, h_dim=128, n_head=4):
3         super().__init__()
4         self.time_embedding = TimeEmbedding(320)
5         self.unet = UNET(h_dim, n_head)
6         self.unet_output = UNETOutputLayer(h_dim, 4)
7
8     @torch.autocast(
9         device_type='cuda', dtype=torch.float16,
10        enabled=True, cache_enabled=True
11    )
12    def forward(self, latent, context, time):
13        time = self.time_embedding(time)
14        output = self.unet(latent, context, time)
15        output = self.unet_output(output)
16        return output
```

Training SD from scratch

❖ Dataset Preparation



```
27 transform = transforms.Compose([
28     transforms.Resize(
29         (WIDTH, HEIGHT),
30         interpolation=transforms.InterpolationMode.BICUBIC
31     ),
32     transforms.ToTensor(),
33     transforms.Normalize(mean=[0.5, 0.5, 0.5], std=[0.5, 0.5, 0.5])
34 ])
35
36 csv_files = ['/content/blobs_crawled_data/metadata.csv']
37 image_folder = '/content/blobs_crawled_data/images'
38
39 train_dataset = EmojiDataset(
40     csv_files=csv_files,
41     image_folder=image_folder,
42     transform=transform
43 )
44 train_dataloader = DataLoader(
45     train_dataset,
46     batch_size=batch_size,
47     shuffle=True,
48     num_workers=2,
49     pin_memory=True,
50     persistent_workers=True
51 )
```



Training SD from scratch

```
1 WIDTH, HEIGHT = 64, 64
2 batch_size = 512
3
4 class EmojiDataset(Dataset):
5     def __init__(self, csv_files, image_folder, transform=None):
6         self.dataframe = pd.concat([pd.read_csv(csv_file) for csv_file in csv_files])
7         self.images_folder = image_folder
8         self.dataframe["image_path"] = self.dataframe["file_name"].str.replace("\\\\", "/")
9         self.image_paths = self.dataframe["image_path"].tolist()
10        self.titles = self.dataframe["prompt"].tolist()
11        self.transform = transform
```

```
13 def __len__(self):
14     return len(self.dataframe)
15
16 def __getitem__(self, idx):
17     image_path = self.images_folder + "/" + self.image_paths[idx]
18     title = self.titles[idx]
19     title = title.replace("'", "").replace('"', "")
20     image = Image.open(image_path).convert('RGB')
21
22     if self.transform:
23         image = self.transform(image)
24
25     return image, title
```

❖ Dataset Preparation

Training SD from scratch

❖ Model Training

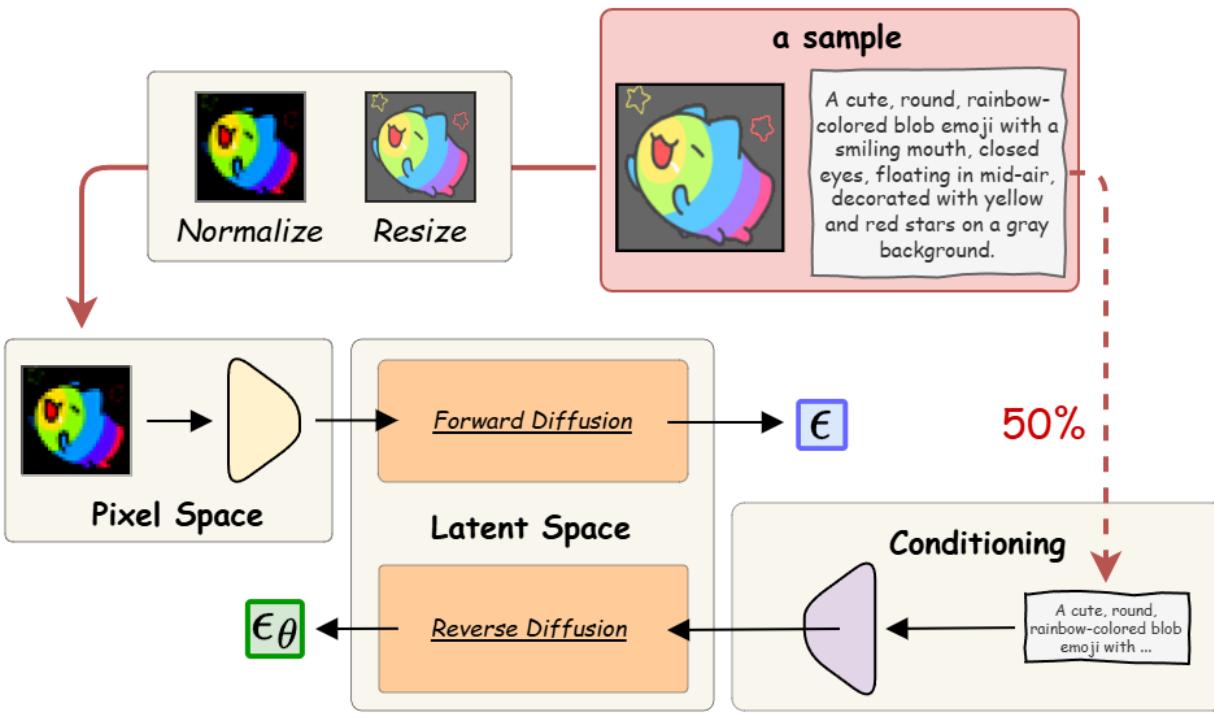
```
1 def train(diffusion, vae, text_encoder, scheduler,
2           optimizer, lr_scheduler, scaler,
3           criterion, dataloader, num_epochs, device="cuda"):
4     losses = []
5     for epoch in range(num_epochs):
6         diffusion.train()
7         epoch_loss = 0.0
8         progress_bar = tqdm(dataloader, desc=f"Epoch {epoch+1}/{num_epochs}", leave=False)
9         for batch_idx, (images, titles) in enumerate(progress_bar):
10            images = images.to(device)
11            image_titles = [f"A photo of {title}" for title in titles]
12            image_titles = [title if random.random() < 0.5 else "" for title in image_titles]
13
14            with torch.no_grad():
15                latents = vae.encode(images).latent_dist.sample() * 0.18215
16                timesteps = torch.randint(
17                    0, scheduler.total_train_timesteps,
18                    (latents.shape[0],), device=device
19                )
20                noisy_latents, noise = scheduler.add_noise(latents, timesteps)
21                time_embeddings = embed_timesteps(timesteps).to(device)
22                text_embeddings = text_encoder(image_titles)
23                noise_pred = diffusion(noisy_latents, text_embeddings, time_embeddings)
24
25                with autocast(device_type='cuda', dtype=torch.float16,
26                              enabled=True, cache_enabled=True):
27                    loss = criterion(noise_pred, noise)
```

```
29     optimizer.zero_grad()
30     scaler.scale(loss).backward()
31     scaler.step(optimizer)
32     scaler.update()
33
34     batch_loss = loss.item()
35     epoch_loss += batch_loss
36
37     progress_bar.set_postfix(loss=f"{batch_loss:.5f}",
38                               lr=f"{optimizer.param_groups[0]['lr']:.6f}")
39
40     lr_scheduler.step()
41     avg_epoch_loss = epoch_loss / len(dataloader)
42     if (epoch + 1) % 10 == 0 or epoch == 0:
43         print(f"Epoch [{epoch+1}/{num_epochs}] - Avg Loss: {avg_epoch_loss:.5f}")
44     losses.append(avg_epoch_loss)
45
46 return losses
```

SKIP

Training SD from scratch

❖ Model Training



$$\text{Loss} = \text{MSE}(\epsilon_\theta, \epsilon)$$

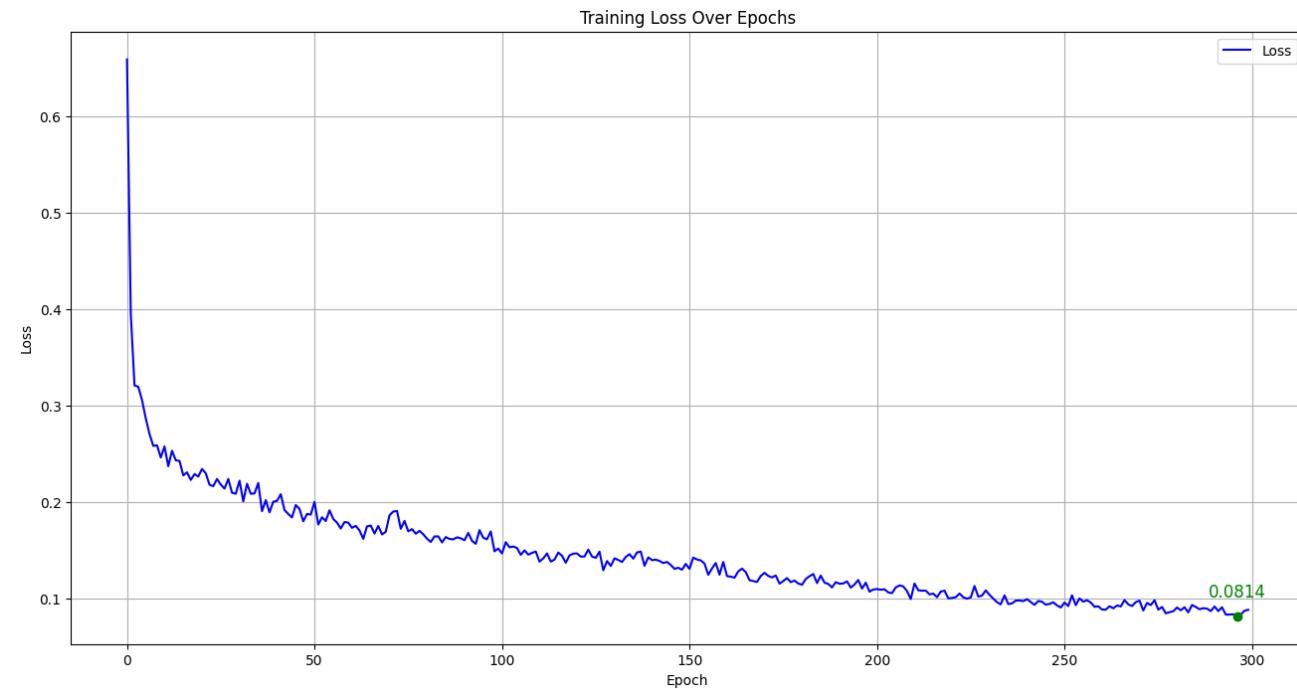
```

1 def train(diffusion, vae, text_encoder, scheduler,
2         optimizer, lr_scheduler, scaler,
3         criterion, dataloader, num_epochs, device="cuda"):
4     losses = []
5     for epoch in range(num_epochs):
6         diffusion.train()
7         epoch_loss = 0.0
8         progress_bar = tqdm(dataloader, desc=f"Epoch {epoch+1}/{num_epochs}", leave=False)
9         for batch_idx, (images, titles) in enumerate(progress_bar):
10            images = images.to(device)
11            image_titles = [f"A photo of {title}" for title in titles]
12            image_titles = [title if random.random() < 0.5 else "" for title in image_titles]
13
14            with torch.no_grad():
15                latents = vae.encode(images).latent_dist.sample() * 0.18215
16                timesteps = torch.randint(
17                    0, scheduler.total_train_timesteps,
18                    (latents.shape[0],), device=device
19                )
20                noisy_latents, noise = scheduler.add_noise(latents, timesteps)
21                time_embeddings = embed_timesteps(timesteps).to(device)
22                text_embeddings = text_encoder(image_titles)
23                noise_pred = diffusion(noisy_latents, text_embeddings, time_embeddings)
24
25                with autocast(device_type='cuda', dtype=torch.float16,
26                              enabled=True, cache_enabled=True):
27                    loss = criterion(noise_pred, noise)

```

Training SD from scratch

❖ Model Training



VAE parameters: 83653863

Diffusion parameters: 72993540

CLIP parameters: 63165952

Total parameters: 219813355

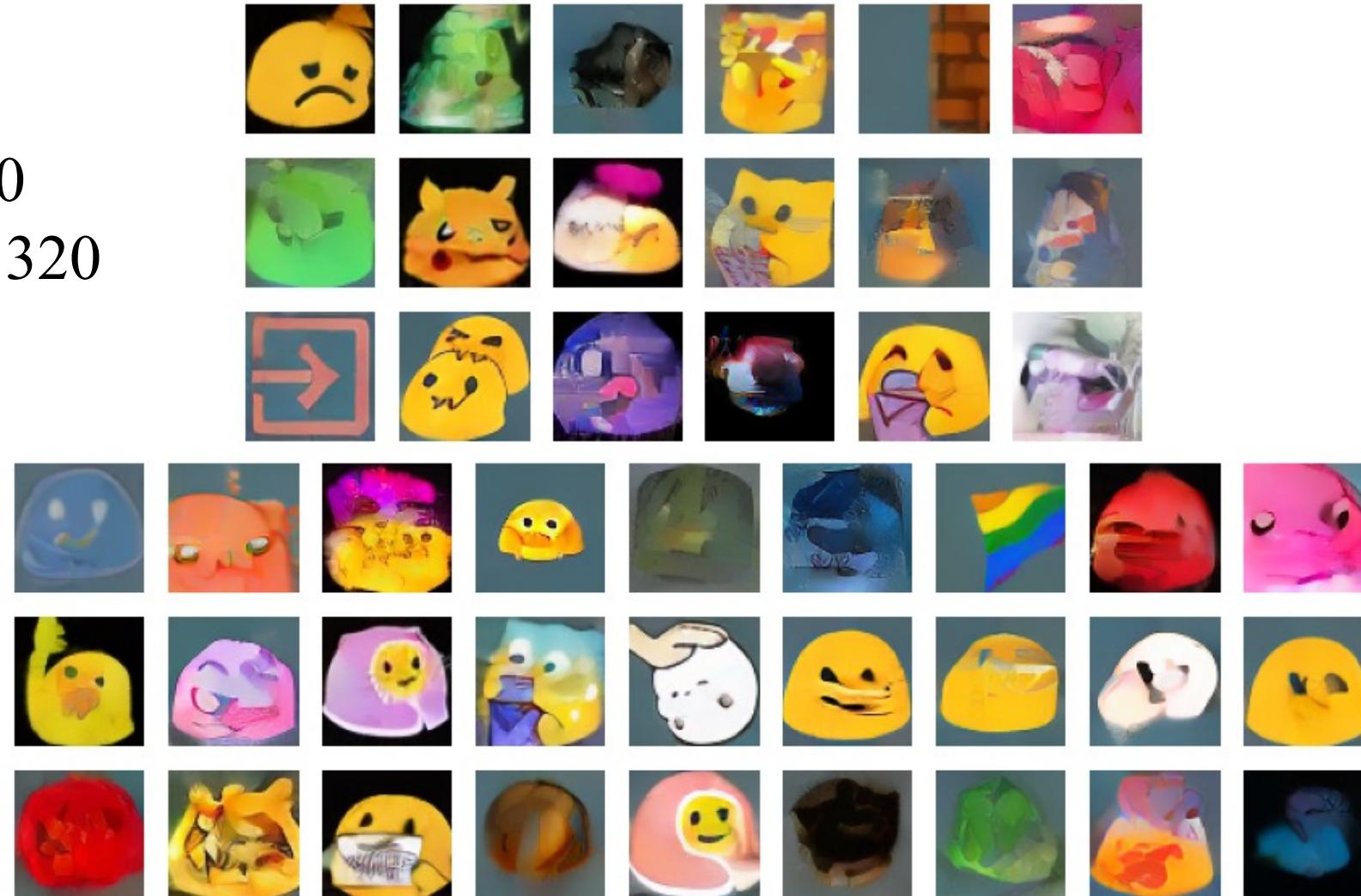
```
1 EPOCHS = 300 # takes more than 3 hours to complete the training
2
3 h_dim = 256
4 n_head = 8
5
6 vae = vae.to(device)
7 diffusion = Diffusion(h_dim, n_head).to(device)
8 clip = CLIPTextEncoder().to(device)
9
10 random_generator = torch.Generator(device='cuda')
11 noise_scheduler = DDPMscheduler(random_generator)
12
13 optimizer = torch.optim.AdamW(diffusion.parameters(), lr=1e-4)
14 criterion = torch.nn.MSELoss()
15
16 lrate_scheduler = lr_scheduler.CosineAnnealingLR(
17     optimizer, T_max=EPOCHS, eta_min=1e-5
18 )
19 scaler = GradScaler()
20
21 def count_parameters(model):
22     return sum(p.numel() for p in model.parameters())
23
24 vae_params = count_parameters(vae)
25 diffusion_params = count_parameters(diffusion)
26 clip_params = count_parameters(clip)
27
28 print(f"VAE parameters: {vae_params}")
29 print(f"Diffusion parameters: {diffusion_params}")
30 print(f"CLIP parameters: {clip_params}")
31 print(f"Total parameters: {vae_params + diffusion_params + clip_params}")
```

Training SD from scratch

❖ Image Generation

Epochs: 300

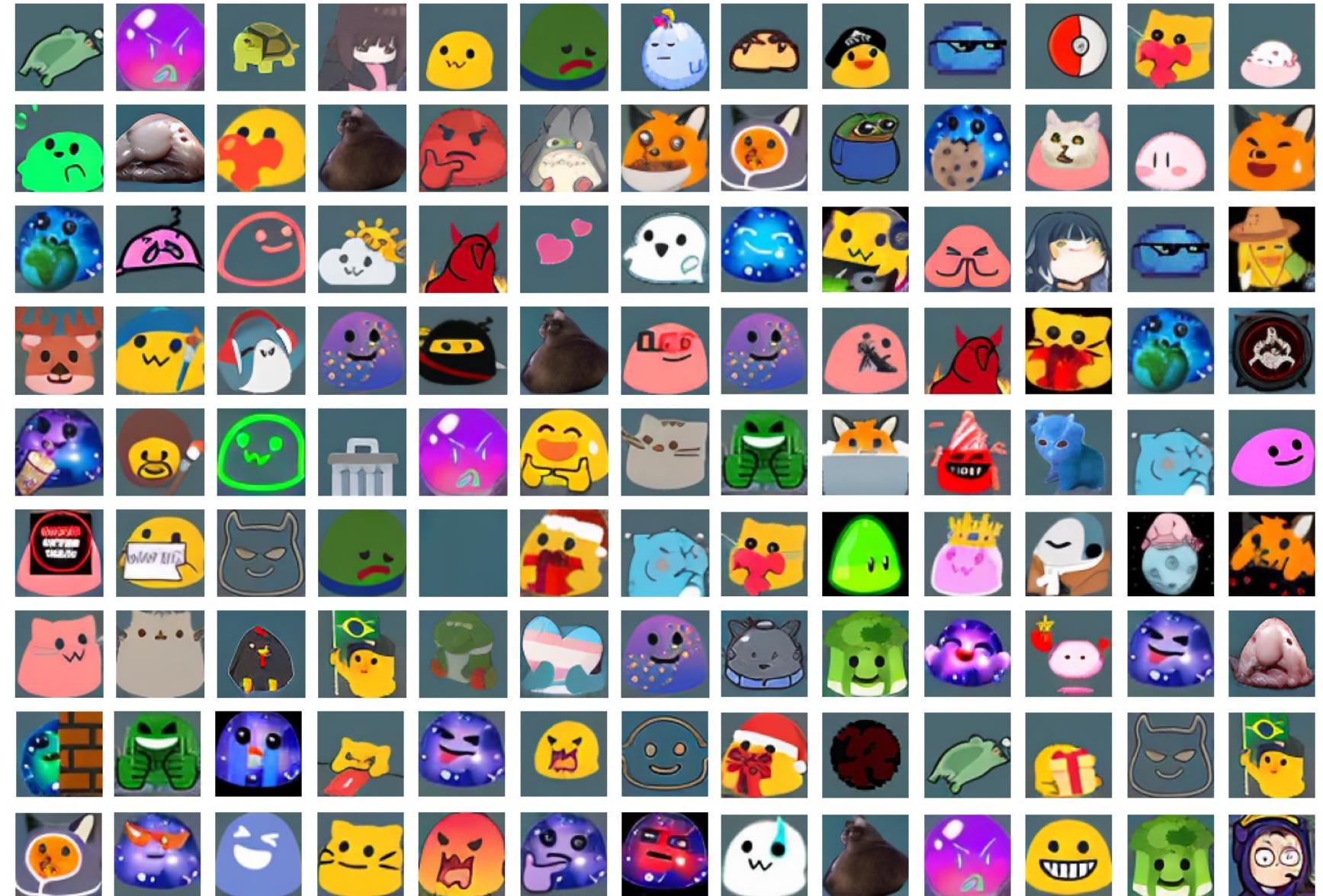
Batch size: 320



Training SD from scratch

❖ Image Generation

Epochs: 5000
Batch size: 640





AI

AI VIET NAM
@aivietnam.edu.vn

QUIZ

Training SD from pre-trained models

Training SD from pre-trained models

❖ Introduction of Stable Diffusion

High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ * Andreas Blattmann¹ * Dominik Lorenz¹ Patrick Esser[✉] Björn Ommer¹

¹Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany [✉]Runway ML

<https://github.com/CompVis/latent-diffusion>

Abstract

By decomposing the image formation process into a sequential application of denoising autoencoders, diffusion models (DMs) achieve state-of-the-art synthesis results on image data and beyond. Additionally, their formulation allows for a guiding mechanism to control the image generation process without retraining. However, since these models typically operate directly in pixel space, optimization of powerful DMs often consumes hundreds of GPU days and inference is expensive due to sequential evaluations. To enable DM training on limited computational resources while retaining their quality and flexibility, we apply them in the latent space of powerful pretrained autoencoders. In contrast to previous work, training diffusion models on such a representation allows for the first time to reach a near-optimal point between complexity reduction and detail preservation, greatly boosting visual fidelity. By introducing cross-attention layers into the model architecture, we turn diffusion models into powerful and flexi-



Figure 1. Boosting the upper bound on achievable quality with less aggressive downsampling. Since diffusion models offer excellent inductive biases for spatial data, we do not need the heavy spatial downsampling of related generative models in latent space, but can still greatly reduce the dimensionality of the data via suitable autoencoding models, see Sec. 3. Images are from the DIV2K [1] validation set, evaluated at 512^2 px. We denote the spatial downsampling factor by f . Reconstruction FIDs [29] and PSNR are calculated on ImageNet-val. [12]; see also Tab. 8.

Stable Diffusion

Stable Diffusion was made possible thanks to a collaboration with [Stability AI](#) and [Runway](#) and builds upon our previous work:

High-Resolution Image Synthesis with Latent Diffusion Models

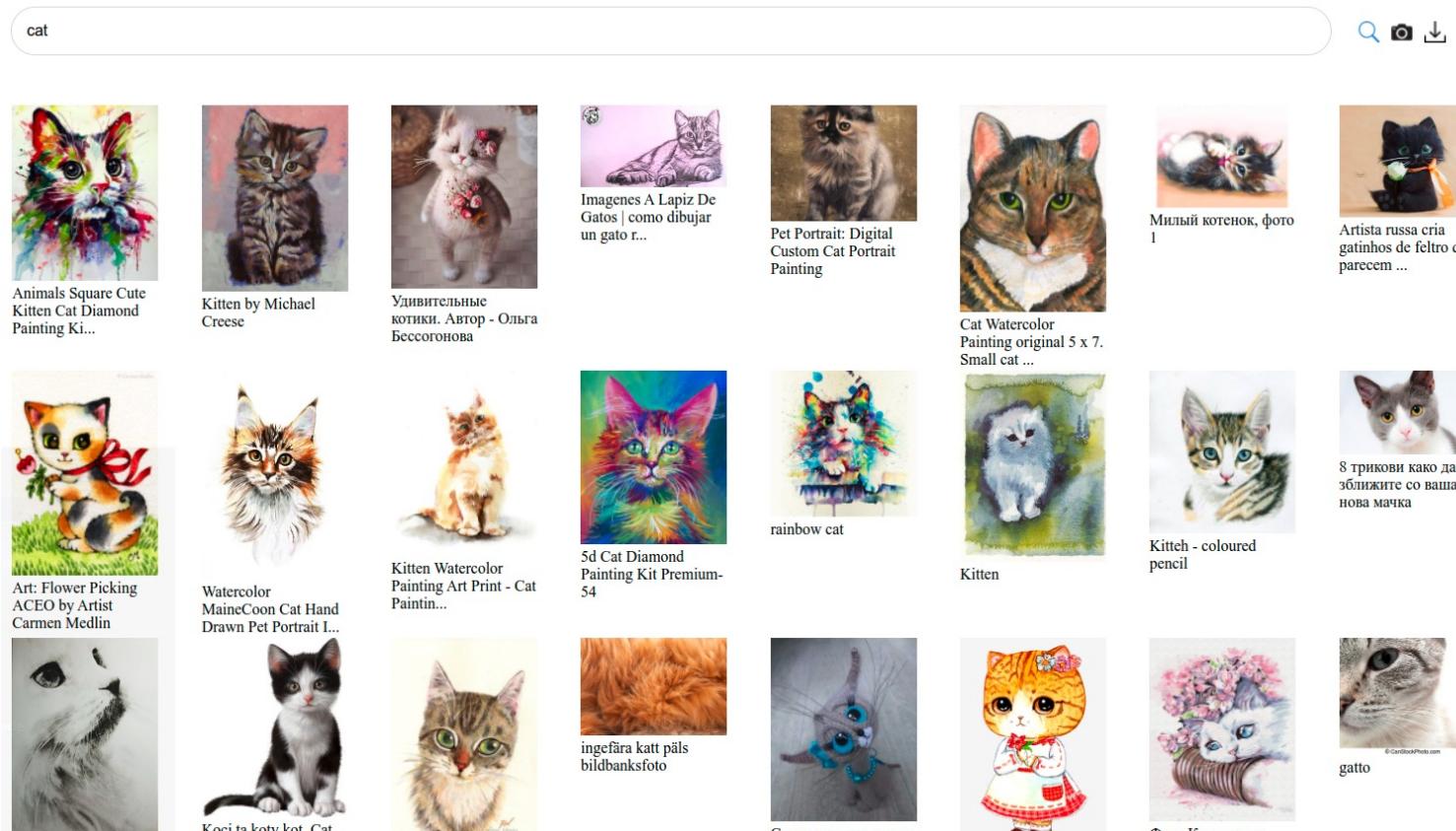
Robin Rombach*, Andreas Blattmann*, Dominik Lorenz, Patrick Esser, Björn Ommer
[CVPR '22 Oral](#) | [GitHub](#) | [arXiv](#) | [Project page](#)



Stable Diffusion is a latent text-to-image diffusion model. Thanks to a generous compute donation from [Stability AI](#) and support from [LAION](#), we were able to train a Latent Diffusion Model on 512x512 images from a subset of the [LAION-5B](#) database. Similar to Google's [Imagen](#), this model uses a frozen CLIP ViT-L/14 text encoder to condition the model on text prompts. With its 860M UNet and 123M text encoder, the model is relatively lightweight and runs on a GPU with at least 10GB VRAM. See [this section](#) below and the [model card](#).

Training SD from pre-trained models

❖ Introduction of Stable Diffusion



Stable Diffusion:

- Text Conditioning:** Incorporates a CLIP-based encoder for stronger text-image alignment.
- Scalability:** Leverages large-scale datasets and extensive training to achieve high-quality image synthesis with fewer computational demands compared to prior latent diffusion variants.

LAION-5: a large, open dataset containing billions of image-text pairs for multi-modal AI research.

Training SD from pre-trained models

❖ Coding: Create PyTorch Datasets

```
1 import os
2 import argparse
3 import math
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import torch
7 import torch.nn as nn
8
9 from torch.utils.data import Dataset, DataLoader
10 from torch.optim.lr_scheduler import LambdaLR
11 from torchvision import transforms
12 from diffusers import StableDiffusionPipeline,
13 from diffusers import UNet2DConditionModel
14 from diffusers import DDPMscheduler
15 from diffusers import AutoencoderKL
16 from transformers import CLIPTextModel, CLIPTokenizer, CLIPFeatureExtractor
17 from huggingface_hub import login, HfApi, create_repo
18 from PIL import Image
19 from tqdm.auto import tqdm
```



Dynamite

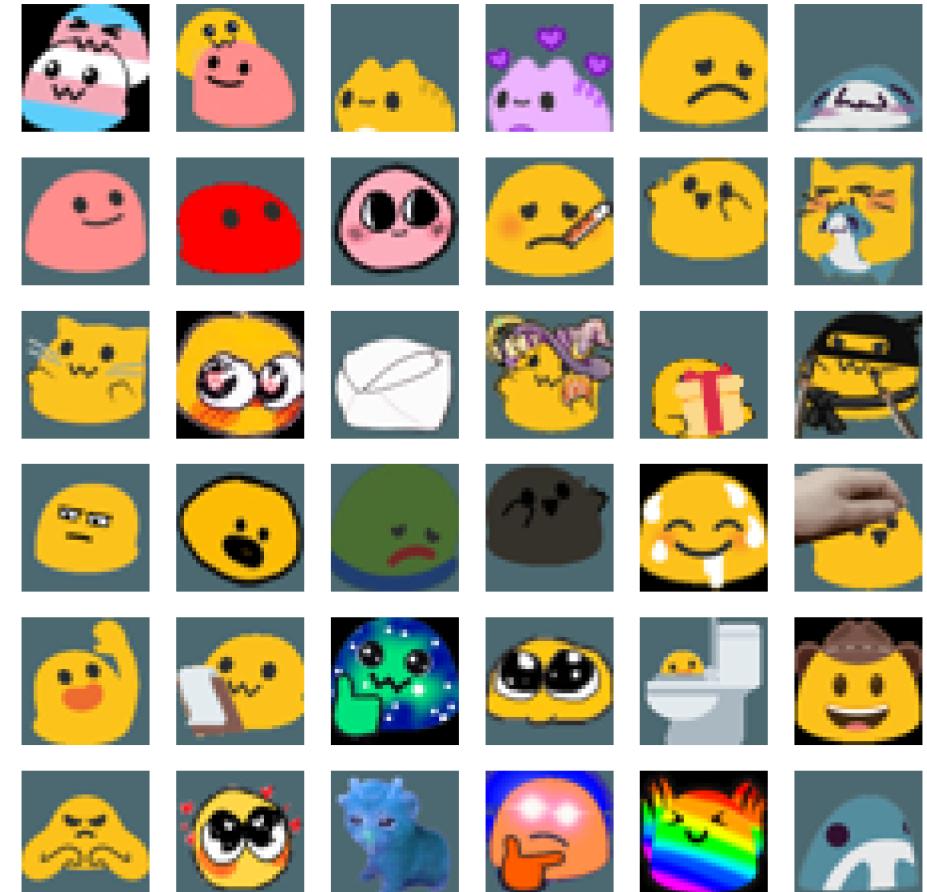


Transformers

Training SD from pre-trained models

❖ Coding: Create PyTorch Datasets

```
1 class EmojiDataset(Dataset):
2     def __init__(self, metadata_df, image_dir, transform=None):
3         self.metadata = metadata_df
4         self.image_dir = image_dir
5         self.transform = transform
6         self.background_color = (255, 255, 255)
7
8     def __len__(self):
9         return len(self.metadata)
10
11    def __getitem__(self, idx):
12        img_path = os.path.join(self.image_dir, self.metadata.iloc[idx]['file_name'])
13        img = Image.open(img_path)
14        img.load()
15        if img.mode == 'RGBA':
16            background = Image.new("RGB", img.size, self.background_color)
17            background.paste(img, mask=img.split()[3])
18            img = background
19        elif 'transparency' in img.info:
20            img = img.convert("RGBA")
21            background = Image.new("RGB", img.size, self.background_color)
22            background.paste(img, mask=img.split()[3])
23            img = background
24        else:
25            img = img.convert("RGB")
26        prompt = self.metadata.iloc[idx]['prompt']
27        if self.transform:
28            img = self.transform(img)
29        return img, prompt
```



Training SD from pre-trained models

❖ Coding: Create PyTorch DataLoader

```
1 def create_dataloaders(data_dir, image_size, batch_size, num_workers=2):
2     metadata = pd.read_csv(os.path.join(data_dir, "metadata.csv"))
3     torch.manual_seed(42)
4     indices = torch.randperm(len(metadata)).tolist()
5     train_size = int(0.9 * len(metadata))
6     test_size = len(metadata) - train_size
7     train_metadata = metadata.iloc[indices[:train_size]].reset_index(drop=True)
8     test_metadata = metadata.iloc[indices[train_size:]].reset_index(drop=True)
9
10    train_transform = transforms.Compose([
11        transforms.Resize((image_size, image_size)),
12        transforms.RandomHorizontalFlip(p=0.5),
13        transforms.ColorJitter(brightness=0.2, contrast=0.2),
14        transforms.ToTensor(),
15        transforms.Normalize([0.5], [0.5])
16    ])
17    test_transform = transforms.Compose([
18        transforms.Resize((image_size, image_size)),
19        transforms.ToTensor(),
20        transforms.Normalize([0.5], [0.5])
21    ])
23    image_dir = os.path.join(data_dir, "images")
24    train_dataset = EmojiDataset(train_metadata, image_dir, train_transform)
25    test_dataset = EmojiDataset(test_metadata, image_dir, test_transform)
26
27    train_loader = DataLoader(
28        train_dataset,
29        batch_size=batch_size,
30        shuffle=True,
31        pin_memory=True,
32        num_workers=num_workers
33    )
34    test_loader = DataLoader(
35        test_dataset,
36        batch_size=batch_size,
37        shuffle=False,
38        pin_memory=True,
39        num_workers=num_workers
40    )
41
42    return train_loader, test_loader, test_dataset
```

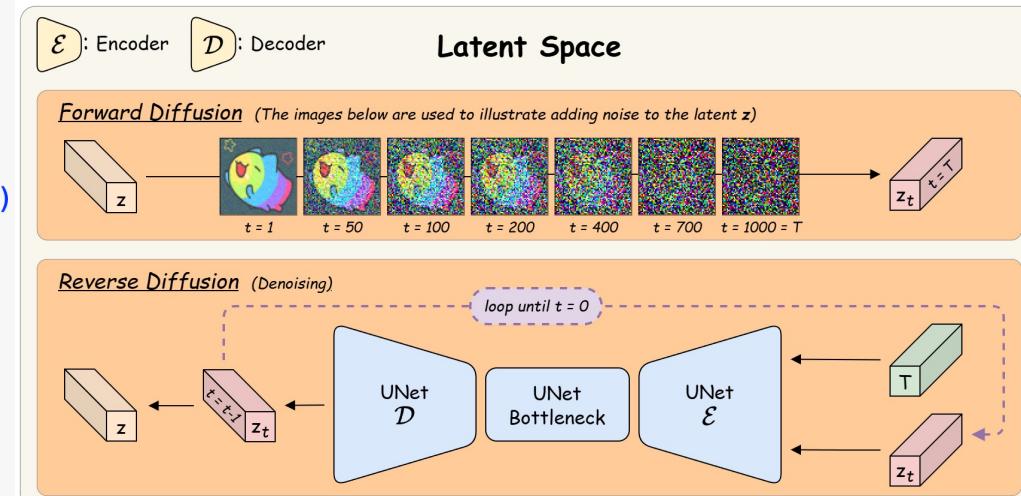
Training SD from pre-trained models

❖ Coding: Create model

```

1 def freeze_unet_layers(unet):
2     for name, param in unet.named_parameters():
3         if "attn2" not in name:
4             param.requires_grad = False
5
6 def build_sd_components(model_id):
7     unet = UNet2DConditionModel.from_pretrained(model_id, subfolder="unet")
8     text_encoder = CLIPTextModel.from_pretrained(model_id, subfolder="text_encoder")
9     tokenizer = CLIPTokenizer.from_pretrained(model_id, subfolder="tokenizer")
10    scheduler = DDPMscheduler.from_pretrained(model_id, subfolder="scheduler")
11    vae = AutoencoderKL.from_pretrained(model_id, subfolder="vae")
12    return unet, text_encoder, tokenizer, scheduler, vae
13
14 def generate_model_name(args, model_id):
15     return (
16         f"sd_bs{args.batch_size}_lr{args.learning_rate}_wd{args.weight_decay}_" +
17         f"_img{args.image_size}_epochs{args.epochs}_{model_id}"
18 )

```





Training SD from pre-trained models

❖ Coding: Create model

CompVis/stable-diffusion-v1-4 like 6.75k Follow CompVis 1.32k

Text-to-Image Diffusers Safetensors StableDiffusionPipeline stable-diffusion stable-diffusion-diffusers arxiv:5 papers License: creativeml-openrail-m

Model card Files and versions xet Community 255 Edit model card

Downloads last month 889,029

Inference Providers NEW HF Inference API Examples

Your sentence here... Compute

View Code Maximize

Model tree for CompVis/stable-diffusion-v1-4 Adapters 525 models Finetunes 1114 models Quantizations 1 model

Spaces using CompVis/stable-diffusion-v1-4 100

Model Details

Stable Diffusion v1-4 Model Card

Stable Diffusion is a latent text-to-image diffusion model capable of generating photo-realistic images given any text input. For more information about how Stable Diffusion functions, please have a look at [😊's Stable Diffusion with 🎨 Diffusers blog](#).

The **Stable-Diffusion-v1-4** checkpoint was initialized with the weights of the [Stable-Diffusion-v1-2](#) checkpoint and subsequently fine-tuned on 225k steps at resolution 512x512 on "laion-aesthetics v2 5+" and 10% dropping of the text-conditioning to improve [classifier-free guidance sampling](#).

This weights here are intended to be used with the [Diffusers](#) library. If you are looking for the weights to be loaded into the CompVis Stable Diffusion codebase, [come here](#)

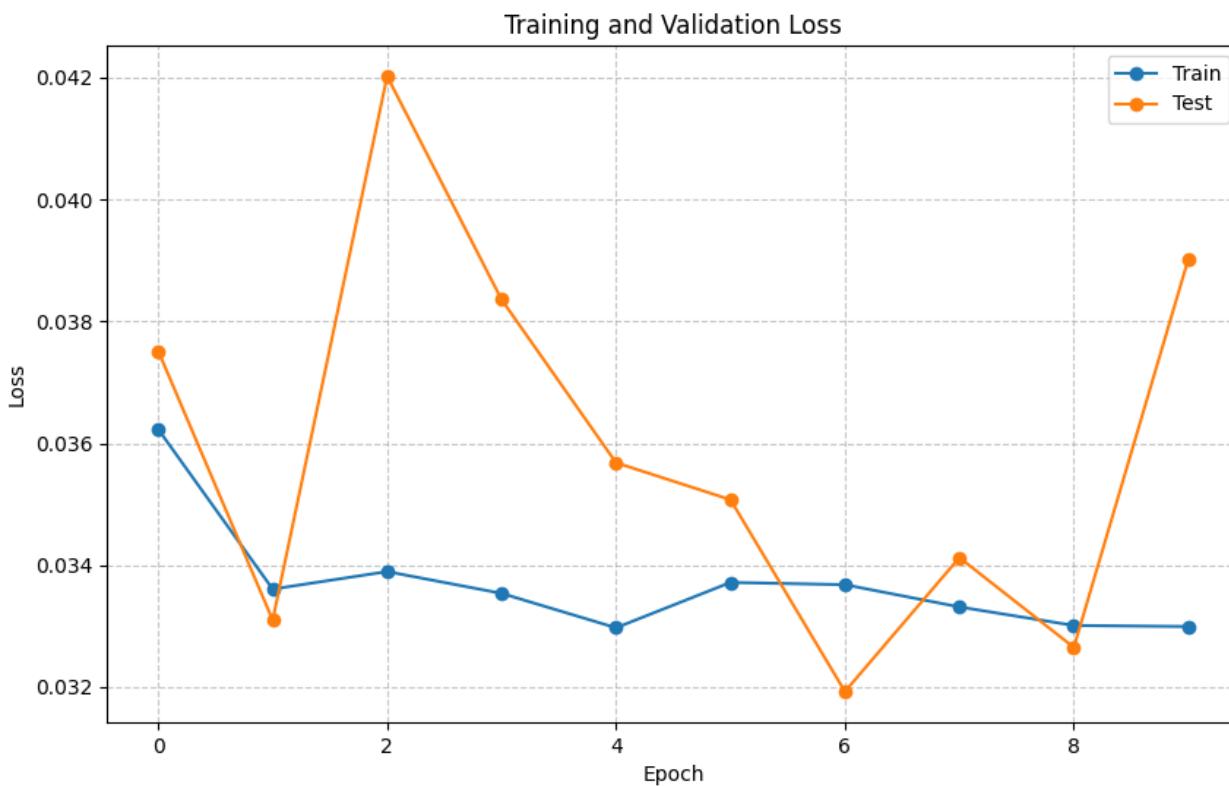
Training SD from pre-trained models

❖ Coding: Create model

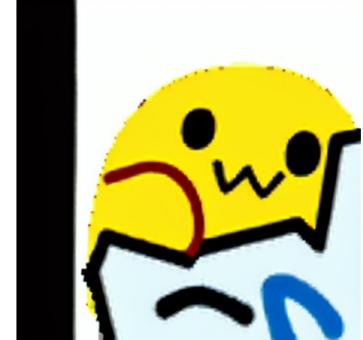
Model	Library	Details
stable-diffusion-v1-1	🤗 Diffusers	237k steps at resolution 256x256 on laion2B-en. 194k steps at resolution 512x512 on laion-high-resolution.
stable-diffusion-v1-2	🤗 Diffusers	v1-1 plus: 515k steps at 512x512 on "laion-improved-aesthetics".
stable-diffusion-v1-3	🤗 Diffusers	v1-2 plus: 195k steps at 512x512 on "laion-improved-aesthetics", with 10% dropping of text-conditioning.
stable-diffusion-v1-4	🤗 Diffusers	v1-2 plus: 225k steps at 512x512 on "laion-aesthetics v2 5+", with 10% dropping of text conditioning.
stable-diffusion-v-1-1-original	CompVis	237k steps at resolution 256x256 on laion2B-en. 194k steps at resolution 512x512 on laion-high-resolution.
stable-diffusion-v-1-2-original	CompVis	v1-1 plus: 515k steps at 512x512 on "laion-improved-aesthetics".
stable-diffusion-v-1-3-original	CompVis	v1-2 plus: 195k steps at 512x512 on "laion-improved-aesthetics", with 10% dropping of text-conditioning.
stable-diffusion-v-1-4-original	CompVis	v1-2 plus: 225k steps at 512x512 on "laion-aesthetics v2 5+", with 10% dropping of text conditioning.

Training SD from pre-trained models

❖ Results



Original: 6019_blobpeek emoji with a round head, small black dot eyes,... Generated: 6019_blobpeek emoji with a round head, small black dot eyes,...



Original: blobcatgay emoji with a round head, rainbow-colored eyes, wi... Generated: blobcatgay emoji with a round head, rainbow-colored eyes, wi...



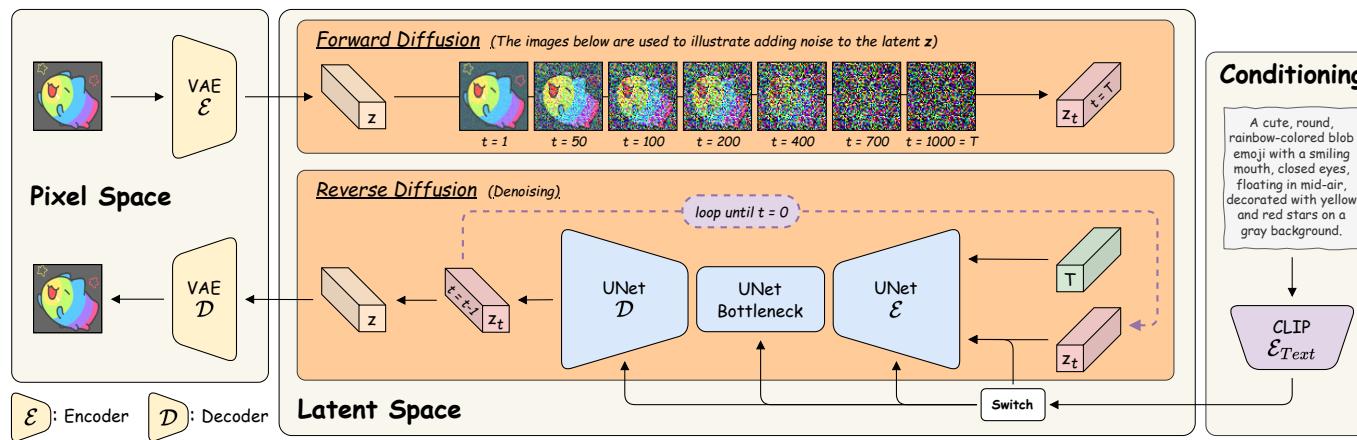
Original: BlobCatBlush2 emoji with a rounded head, closed eyes, slight... Generated: BlobCatBlush2 emoji with a rounded head, closed eyes, slight...



Summarization and Q&A

Summarization and Q&A

❖ Content



In this section, we have discussed about:

- The task of Image Generation (IG) and Stable Diffusion (SD).
- A Python program to crawl Discord Emojis on website.
- Implement and training:
 - A SD model from scratch using PyTorch and Diffusers.
 - A SD model from a pre-trained model using Diffusers.

Summarization and Q&A

