



CLEVELAND STATE UNIVERSITY

CIS 660 – DATA MINING

PROJECT PRESENTATION

GROUP MEMBERS:
RAHUL PADWANI - 2873059
JAYKUMAR MISTRY - 2875386

PROFESSOR: - DR. SUNNIE CHUNG

TEXT-TO-IMAGE SYNTHESIS WITH DEEP NEURAL NETWORKS



INTRODUCTION

IN OUR PROJECT, WE'RE TURNING THE VISION OF WORDS PAINTING PICTURES INTO REALITY. USING A COLLECTION OF **31,000 IMAGES** LINKED TO DESCRIPTIVE SENTENCES, WE'RE TRAINING A COMPUTER TO BRIDGE THE GAP BETWEEN WORDS AND VISUALS. THIS GOES BEYOND ART CREATION – IT'S ABOUT TRANSFORMING LANGUAGE INTO VIVID IMAGES THROUGH A SOPHISTICATED PROGRAM, LIKE A **PIXEL-BY-PIXEL CONVERSATION** WITH A PAINTER. THE TECHNOLOGY'S POTENTIAL SPANS GRAPHIC DESIGN CREATIVITY, ACCESSIBILITY REVOLUTION, AND THE TRANSFORMATION OF EDUCATION AND ENTERTAINMENT INTO INTERACTIVE **VISUAL EXPERIENCES**. OUR PROJECT IS MORE THAN TECHNOLOGICAL ADVANCEMENT; IT'S ABOUT DEMOCRATIZING THE POWER OF VISUAL STORYTELLING, MAKING IT ACCESSIBLE TO EVERYONE.

GOAL OF THE PROJECT

- **PRIMARY GOAL:** DEVELOP AN ADVANCED MACHINE LEARNING MODEL CAPABLE OF ACCURATELY GENERATING IMAGES FROM TEXTUAL DESCRIPTIONS.
- **INTEGRATION OF TECHNIQUES:** COMBINE CUTTING-EDGE TECHNIQUES IN COMPUTER VISION, NATURAL LANGUAGE PROCESSING AND STABLE DIFFUSION MODEL FOR A ROBUST TEXT-TO-IMAGE SYNTHESIS SYSTEM.
- **OVERALL OBJECTIVE:** CREATE A VERSATILE AND EFFECTIVE TOOL THAT SEAMLESSLY TRANSLATES DESCRIPTIVE LANGUAGE INTO VISUALLY COMPELLING IMAGES.

SYSTEMS AND TOOLS USED

DEVELOPMENT ENVIRONMENT: **GOOGLE COLAB** AND **JUPYTER NOTEBOOK**

PROGRAMMING LANGUAGE: **PYTHON**

LIBRARIES: **TENSORFLOW**, **PyTorch**, **NLTK**, **PIL**, **OPENCV...**

MACHINE LEARNING ALGORITHMS: **CNNs**, **GANs**, **STABLE DIFFUSION**

DATA COLLECTION

DATA COLLECTION

DATASET UTILIZATION: UTILIZE THE EXTENSIVE FLICKR30K DATASET TO ENABLE THE MODEL TO LEARN COMPLEX RELATIONSHIPS BETWEEN TEXT AND VISUAL ELEMENTS.

DATASET SOURCE:

KAGGLE -[HTTPS://WWW.KAGGLE.COM/DATASETS/ADITYAJN105/Flickr30k/DATA](https://www.kaggle.com/datasets/adityajn105/flickr30k/data)

DATASET OVERVIEW:

ROBUST REPRESENTATION WITH APPROXIMATELY **155,000 DATA POINTS**.

EACH IMAGE-DESCRIPTION PAIR IS TREATED AS AN INDIVIDUAL DATA POINT.

COMPRISSES 31,000 IMAGES, WITH **EACH IMAGE PAIRED WITH FIVE DISTINCT TEXT DESCRIPTIONS**.

DATASET SIZE:

AROUND 155,000 TEXT-IMAGE PAIRS ($31,000 \text{ IMAGES} \times 5 \text{ CAPTIONS/IMAGE}$).

SIZE OF THE DATASET IS APPROXIMATELY 8.86 GB.

STRUCTURE:

THE DATASET INCLUDES AN **IMAGES** FOLDER CONTAINING **FLICKER30K_IMAGES**.

THE "**CAPTION.TXT**" FILE HOLDS FIVE DIFFERENT TEXT DESCRIPTIONS PER IMAGE, PROVIDING DIVERSE PERSPECTIVES ON THE CONTENTS AND CONTEXT OF THE IMAGES.

Flick 30k Dataset

Flick 30k Dataset for Image Captioning

Data Card Code (5) Discussion (0)

About Dataset

Context
A new benchmark collection for sentence-based image description and search, consisting of 30,000 images that are each paired with five different captions which provide clear descriptions of the salient entities and events. ... The images were chosen from six different Flickr groups, and tend not to contain any well-known people or locations, but were manually selected to depict a variety of scenes and situations

Content
What's inside is more than just rows and columns. Make it easy for others to get started by describing how you acquired the data and what time period it represents, too.

Acknowledgements
We wouldn't be here without the help of others. If you owe any attributions or thanks, include them here along with any citations of past research.

Usability 6.25

License CC0: Public Domain

Expected update frequency Not specified

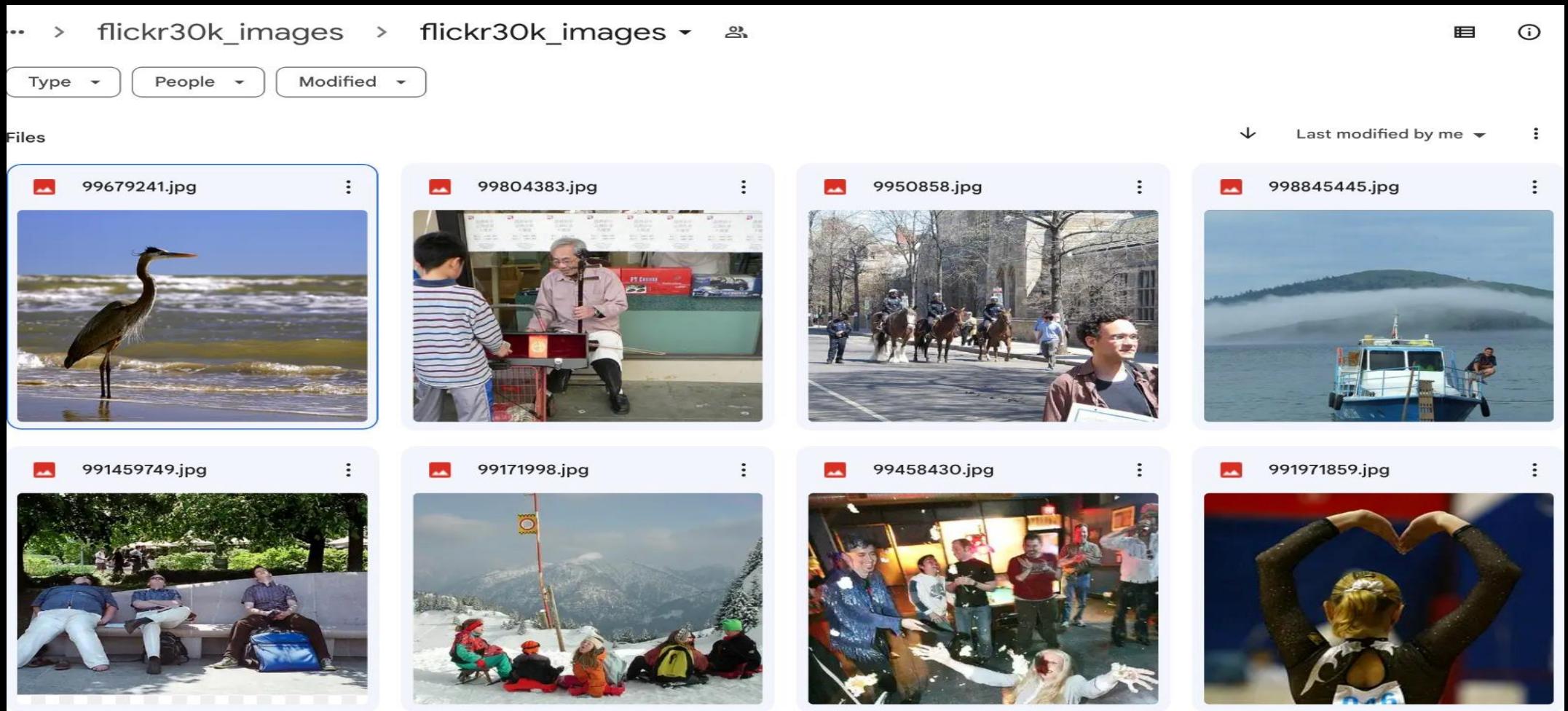
Tags



CAPTION.TXT

Results			
image_name comment_number comment			
1000092795.jpg 0 Two young guys with shaggy hair look at their hands while hanging out in the yard .			
1000092795.jpg 1 Two young		White males are outside near many bushes .	
1000092795.jpg 2 Two men in green shirts are standing in a yard .			
1000092795.jpg 3 A man in a blue shirt standing in a garden .			
1000092795.jpg 4 Two friends enjoy time spent together .			
10002456.jpg 0 Several men in hard hats are operating a giant pulley system .			
10002456.jpg 1 Workers look down from up above on a piece of equipment .			
10002456.jpg 2 Two men working on a machine wearing hard hats .			
10002456.jpg 3 Four men on top of a tall structure .			
10002456.jpg 4 Three men on a large rig .			
1000268201.jpg 0 A child in a pink dress is climbing up a set of stairs in an entry way .			
1000268201.jpg 1 A little girl in a pink dress going into a wooden cabin .			
1000268201.jpg 2 A little girl climbing the stairs to her playhouse .			
1000268201.jpg 3 A little girl climbing into a wooden playhouse			
1000268201.jpg 4 A girl going into a wooden building .			
1000344755.jpg 0 Someone in a blue shirt and hat is standing on stair and leaning against a window .			
1000344755.jpg 1 A man in a blue shirt is standing on a ladder cleaning a window .			
1000344755.jpg 2 A man on a ladder cleans the window of a tall building .			
1000344755.jpg 3 man in blue shirt and jeans on ladder cleaning windows			
1000344755.jpg 4 a man on a ladder cleans a window			
1000366164.jpg 0 Two men	one in a gray shirt	one in a black shirt	standing near a stove .
1000366164.jpg 1 Two guy cooking and joking around with the camera .			
1000366164.jpg 2 Two men in a kitchen cooking food on a stove .			
1000366164.jpg 3 Two men are at the stove preparing food .			
1000366164.jpg 4 Two men are cooking a meal .			

FLICKER30K_IMAGES FOLDER



DATA PREPROCESSING

DATA PREPROCESSING

Part – 1 Image Preprocessing

Part – 2 Text Preprocessing

Part – 3 Combining Text and image in one result

PART – 1 IMAGE PREPROCESSING



```
import os
from PIL import Image
from tqdm.notebook import tqdm
import numpy as np

# Directory where your original images are stored
image_dir = 'flickr30k_images/flickr30k_images'

# Directory where you want to save the preprocessed images
processed_image_dir = 'flickr30k_images/processed_images'
os.makedirs(processed_image_dir, exist_ok=True)

# List to keep track of processed files
processed_files = []
not_processed_files = []

# Function to preprocess images
def transform(image):
    return image
```

Path Description

```
# Optionally, print the filenames of images not processed
if not_processed_files:
    print("Images not processed:")
    for filename in not_processed_files:
        print(filename)

num_images_to_print = 5 # Adjust as needed
for filename in processed_files[:num_images_to_print]:
    npy_filename = os.path.splitext(filename)[0] + '.npy'
    load_path = os.path.join(processed_image_dir, npy_filename)
    image_array = np.load(load_path)
    print(f"Array for {filename}:")
    print(image_array)
```

Printing the image array.

PART – 1 IMAGE PREPROCESSING

```
# Process and save images as NumPy arrays
for filename in tqdm(os.listdir(image_dir)):
    if filename.lower().endswith('.jpg'): # Check if the file is an image
        try:
            file_path = os.path.join(image_dir, filename)
            with Image.open(file_path) as img:
                transformed_img = transform(img)
                npy_filename = os.path.splitext(filename)[0] + '.npy'
                save_path = os.path.join(processed_image_dir, npy_filename)
                np.save(save_path, transformed_img)
                processed_files.append(filename) # Add to the list of processed files
        except Exception as e:
            print(f"Error processing {filename}: {e}")
            not_processed_files.append(filename) # Add to the list if not processed

# Print the results
print(f"Total number of images processed: {len(processed_files)}")
print(f"Total number of images not processed: {len(not_processed_files)}")
```

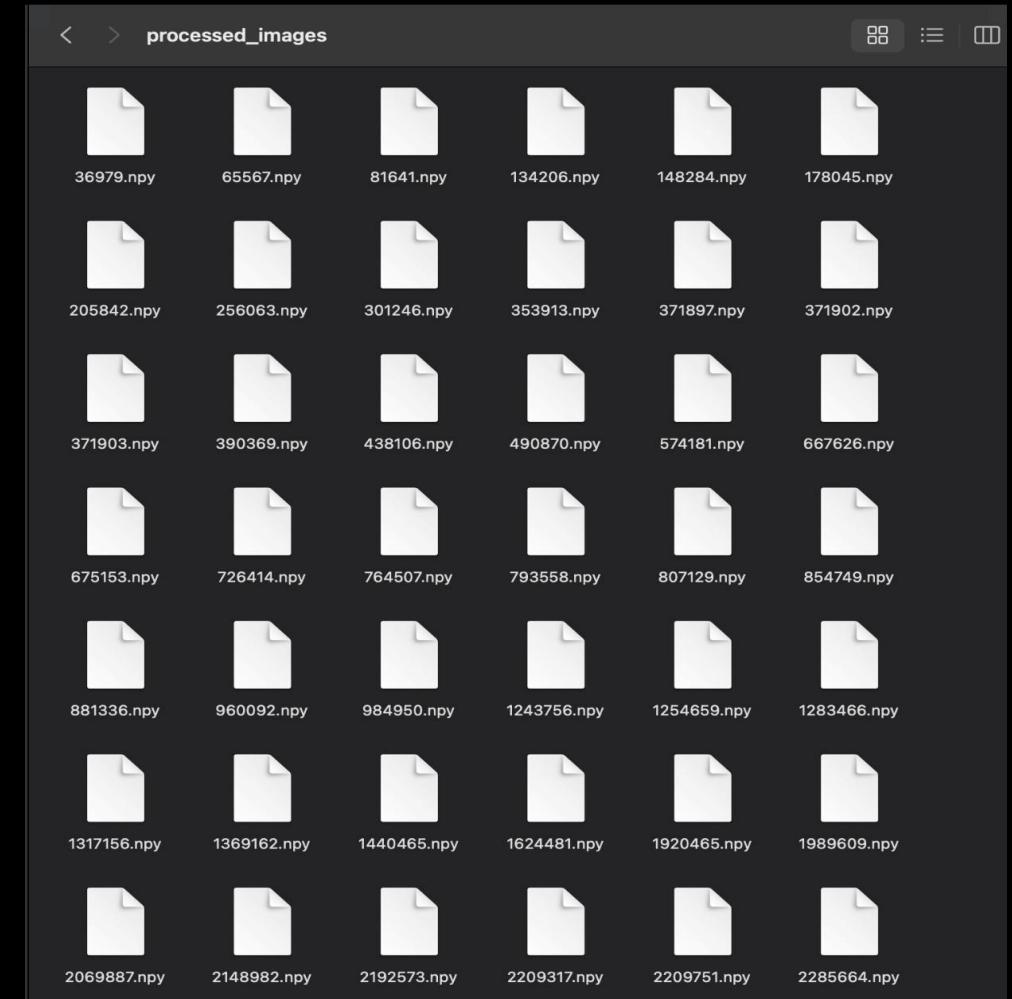
- Checking the image (jpg)
- Resizing the Image.
- Normalization.
- Noise Reduction.
- Splitting extension and adding with the NumPy array
- Printing total numbers of image processed and not processed

PART – 1 IMAGE PREPROCESSING

```
100% 31785/31785 [02:06<00:00, 258.67it/s]

Error processing 4904808403.jpg: cannot identify image file 'flickr30k_images/flickr30k_images/4904808403.jpg'
Total number of images processed: 31782
Total number of images not processed: 1
Images not processed:
4904808403.jpg
Array for 2609797461.jpg:
[[[103 123 148]
 [105 125 150]
 [109 126 152]
 ...
 [198 203 206]
 [196 201 204]
 [196 201 204]]
 [[103 123 148]
 [105 125 150]
 [109 126 152]]]
```

Preprocessed Image



PART – 2 TEXT PREPROCESSING

```
[ ] import pandas as pd
import re

# Path to your original captions file
captions_file_path = 'flickr30k_images/results.csv'

# Read the CSV file
captions_df = pd.read_csv(captions_file_path, delimiter='|')

# Rename the columns to remove leading/trailing spaces
captions_df.columns = captions_df.columns.str.strip()

# Function to clean text
def clean_text(text):
    text = text.lower() # Convert text to lowercase
    text = re.sub(r'^\w\s]', '', text) # Remove punctuation and numbers
    text = re.sub(r'\d+', '', text) # Remove numbers
    text = re.sub(r'\s+', ' ', text).strip() # Remove extra spaces
    return text

# Apply the cleaning function to the comment column
captions_df['comment'] = captions_df['comment'].astype(str).apply(clean_text)
```

captions_df			
	image_name	comment_number	comment
0	1000092795.jpg	0	two young guys with shaggy hair look at their ...
1	1000092795.jpg	1	two young white males are outside near many bu...
2	1000092795.jpg	2	two men in green shirts are standing in a yard
3	1000092795.jpg	3	a man in a blue shirt standing in a garden
4	1000092795.jpg	4	two friends enjoy time spent together
...
158910	998845445.jpg	0	a man in shorts and a hawaiian shirt leans ove...
158911	998845445.jpg	1	a young man hanging over the side of a boat wh...
158912	998845445.jpg	2	a man is leaning off of the side of a blue and...
158913	998845445.jpg	3	a man riding a small boat in a harbor with fog...
158914	998845445.jpg	4	a man on a moored blue and white boat with hil...
158915 rows × 3 columns			

- Separating columns with ‘|’
- Removing punctuation and number
- Removing extra spaces
- Output of the cleaned data

PART – 2 TEXT PREPROCESSING

TOKENIZER AND VECTORIZATION ON TEXT DATA

```
▶ from tensorflow.keras.preprocessing.text import Tokenizer  
from tensorflow.keras.preprocessing.sequence import pad_sequences  
  
# Tokenization and vectorization  
MAX_NUM_WORDS = 10000  
tokenizer = Tokenizer(num_words=MAX_NUM_WORDS)  
tokenizer.fit_on_texts(captions_df['comment'])  
sequences = tokenizer.texts_to_sequences(captions_df['comment'])  
word_index = tokenizer.word_index  
data = pad_sequences(sequences, maxlen=100)
```

Converting Comments in matrices for machine
To understand.

	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8
Term(s) 1	10	0	1	0	0	0	0	2
Term(s) 2	0	2	0	0	0	18	0	2
Term(s) 3	0	0	0	0	0	0	0	2
Term(s) 4	6	0	0	4	6	0	0	0
Term(s) 5	0	0	0	0	0	0	0	2
Term(s) 6	0	0	1	0	0	1	0	0
Term(s) 7	0	1	8	0	0	0	0	0
Term(s) 8	0	0	0	0	0	3	0	0

Word Vector (Passage Vector)

Document Vector

PART – 3 COMBINING TEXT AND IMAGE IN ARRAY

combined_dataset	
image_filename	text_data
516279719	[0 11 18 321 9 2098 107 183 15 60 155 22 320 3 485]
7610186624	[0 192 1416]
2735799584	[0 11 28 2 45 252 12 30 2 1 485]
5511458367	[0 1 6 2 1 24 20 30 2 1 705]
2044546977	[0 9961 134]
3726980861	[0 3831 2622]

```
▶ image_filenames = [os.path.splitext(filename)[0] for filename in os.listdir(processed_image_dir) if filename.endswith('.npy')]

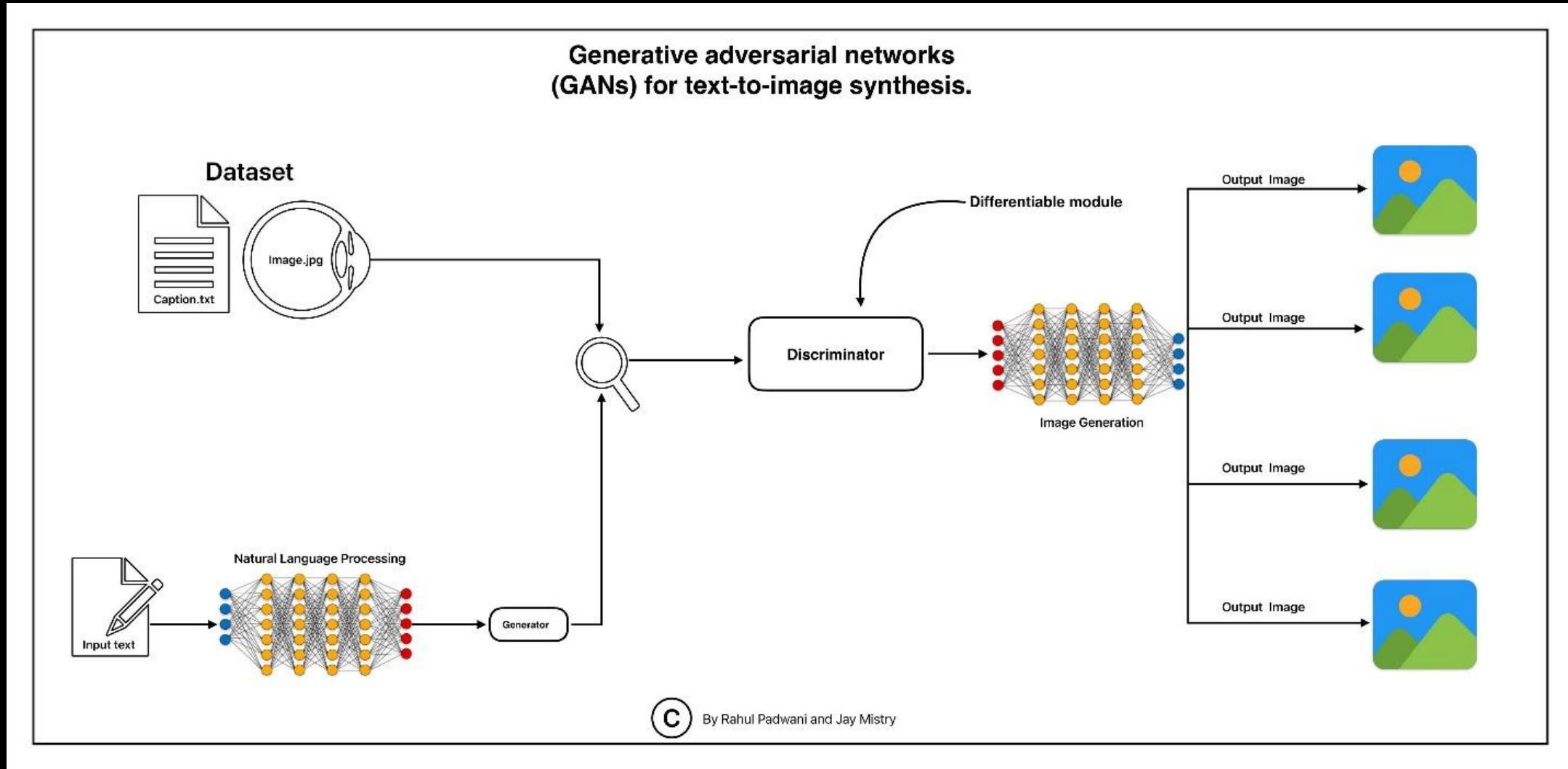
# Create DataFrame
export_df = pd.DataFrame({
    'image_filename': image_filenames,
    'text_data': [text for _, text in combined_dataset]
})

# Export to CSV
export_df.to_csv('flickr30k_images/combined_dataset.csv', index=False)
```

As per the image name the text data is mapped and converted into standard format.

GAN IMPLEMENTATION

GAN IMPLEMENTATION



GAN IMPLEMENTATION – GENERATOR

```
def make_generator_model():
    model = tf.keras.Sequential()
    model.add(layers.Dense(8*8*256, use_bias=False, input_shape=(100,)))
    model.add(layers.BatchNormalization())
    model.add(layers.LeakyReLU())

    model.add(layers.Reshape((8, 8, 256))) # Reshape to a 3D tensor
    assert model.output_shape == (None, 8, 8, 256)

    # Upsample to 16x16
    model.add(layers.Conv2DTranspose(128, (5, 5), strides=(2, 2), padding='same', use_bias=False))
    assert model.output_shape == (None, 16, 16, 128)
    model.add(layers.BatchNormalization())
    model.add(layers.LeakyReLU())

    # Upsample to 32x32
    model.add(layers.Conv2DTranspose(64, (5, 5), strides=(2, 2), padding='same', use_bias=False))
    assert model.output_shape == (None, 32, 32, 64)
    model.add(layers.BatchNormalization())
    model.add(layers.LeakyReLU())

    # Upsample to 64x64
    model.add(layers.Conv2DTranspose(3, (5, 5), strides=(2, 2), padding='same', use_bias=False, activation='tanh'))
    assert model.output_shape == (None, 64, 64, 3)

    return model
```

The generator model consists of a dense layer with 16,384 units (1,638,400 parameters), batch normalization (65,536 parameters), Leaky ReLU activation, a reshape layer for a 4D tensor, multiple Conv2D Transpose layers for upsampling, totaling 2,733,504 parameters (2,700,352 trainable).

Model: "sequential_4"			
Layer (type)	Output Shape	Param #	
dense_4 (Dense)	(None, 16384)	1638400	
batch_normalization_6 (BatchNormalization)	(None, 16384)	65536	
leaky_re_lu_10 (LeakyReLU)	(None, 16384)	0	
reshape_2 (Reshape)	(None, 8, 8, 256)	0	
conv2d_transpose_6 (Conv2DTranspose)	(None, 16, 16, 128)	819200	
batch_normalization_7 (BatchNormalization)	(None, 16, 16, 128)	512	
leaky_re_lu_11 (LeakyReLU)	(None, 16, 16, 128)	0	
conv2d_transpose_7 (Conv2DTranspose)	(None, 32, 32, 64)	204800	
batch_normalization_8 (BatchNormalization)	(None, 32, 32, 64)	256	
leaky_re_lu_12 (LeakyReLU)	(None, 32, 32, 64)	0	
conv2d_transpose_8 (Conv2DTranspose)	(None, 64, 64, 3)	4800	
=====			
Total params: 2733504 (10.43 MB)			
Trainable params: 2700352 (10.30 MB)			
Non-trainable params: 33152 (129.50 KB)			
None			

GAN IMPLEMENTATION – DISCRIMINATOR

```
def make_discriminator_model():
    model = tf.keras.Sequential()
    model.add(layers.Conv2D(64, (5, 5), strides=(2, 2), padding='same', input_shape=[64, 64, 3]))
    model.add(layers.LeakyReLU())
    model.add(layers.Dropout(0.3))

    model.add(layers.Conv2D(128, (5, 5), strides=(2, 2), padding='same'))
    model.add(layers.LeakyReLU())
    model.add(layers.Dropout(0.3))

    model.add(layers.Flatten())
    model.add(layers.Dense(1))

    return model

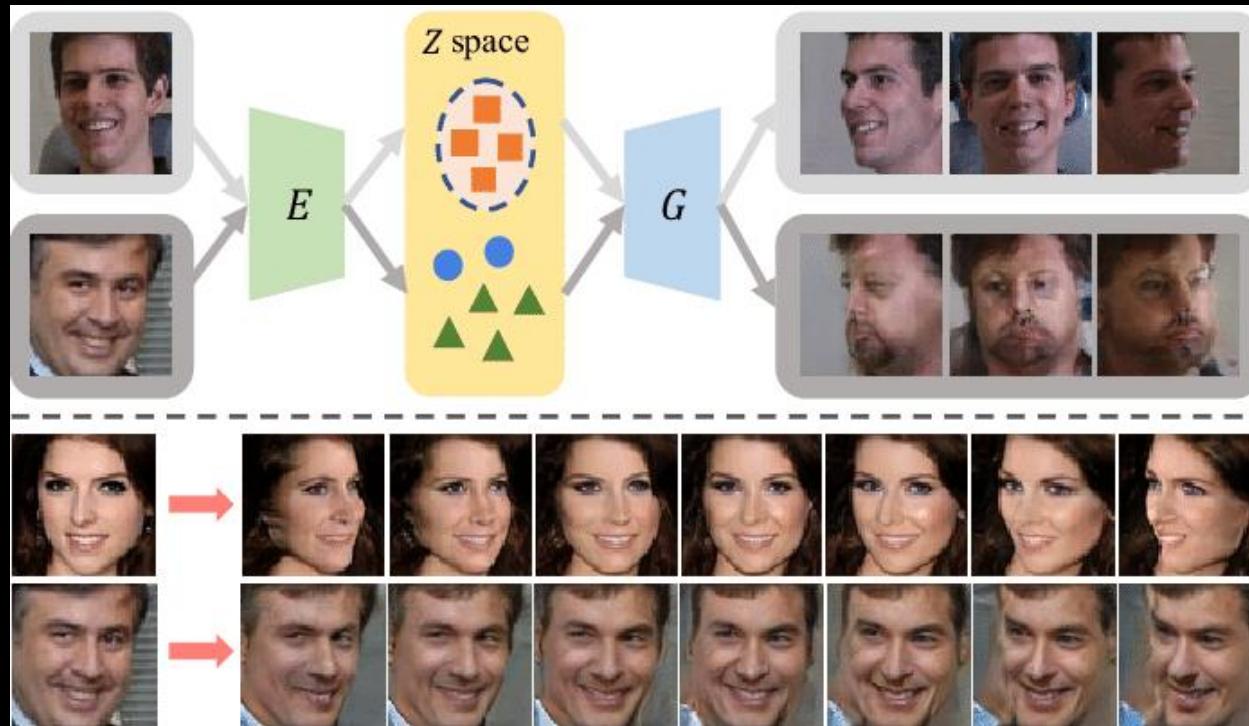
generator = make_generator_model()
discriminator = make_discriminator_model()

# Print model summaries
print(generator.summary())
print(discriminator.summary())
```

The discriminator model features convolutional layers (Conv2D) for downsampling (209,792 parameters), Leaky ReLU activation, dropout for overfitting prevention, a flattening layer, and a dense layer for output, totaling 242,561 trainable parameters.

Model: "sequential_5"		
Layer (type)	Output Shape	Param #
conv2d_4 (Conv2D)	(None, 32, 32, 64)	4864
leaky_re_lu_13 (LeakyReLU)	(None, 32, 32, 64)	0
dropout_4 (Dropout)	(None, 32, 32, 64)	0
conv2d_5 (Conv2D)	(None, 16, 16, 128)	204928
leaky_re_lu_14 (LeakyReLU)	(None, 16, 16, 128)	0
dropout_5 (Dropout)	(None, 16, 16, 128)	0
flatten_2 (Flatten)	(None, 32768)	0
dense_5 (Dense)	(None, 1)	32769
Total params: 242561 (947.50 KB)		
Trainable params: 242561 (947.50 KB)		
Non-trainable params: 0 (0.00 Byte)		
None		

GAN IMPLEMENTATION – LIMITATION



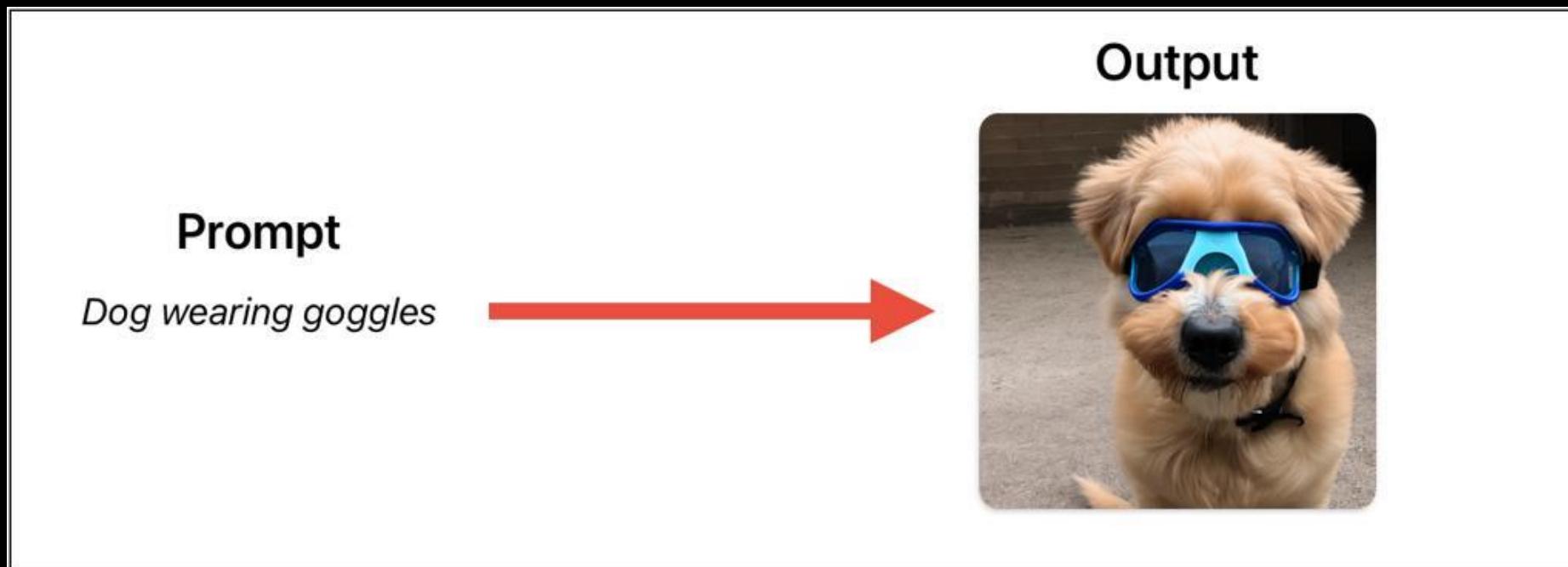
In contrast, the promising results of GANs have been revealed to be mostly confined to data with comparably limited variability as their adversarial learning procedure does not easily scale to modeling complex, multi-modal distributions.

STABLE DIFFUSION

STABLE DIFFUSION

WHAT IS STABLE DIFFUSION?

STABLE DIFFUSION IS AN ADVANCED DEEP LEARNING MODEL USED FOR GENERATING HIGH-QUALITY IMAGES FROM TEXTUAL DESCRIPTIONS. IT OPERATES ON THE PRINCIPLES OF DIFFUSION MODELS, A TYPE OF GENERATIVE MODEL THAT GRADUALLY TRANSFORMS A RANDOM PATTERN OF PIXELS INTO A COHERENT IMAGE.



WHAT IS STABLE DIFFUSION?

TEXT-TO-IMAGE SYNTHESIS: STABLE DIFFUSION CAN CREATE DETAILED AND REALISTIC IMAGES BASED ON TEXTUAL PROMPTS, MAKING IT A POWERFUL TOOL FOR ARTISTS, DESIGNERS, AND CONTENT CREATORS.

DIFFUSION PROCESS: IT USES A PROCESS WHERE IT STARTS WITH A RANDOM NOISE IMAGE AND GRADUALLY REFINES IT TO MATCH THE DESIRED OUTPUT, GUIDED BY THE TEXTUAL INPUT.

VERSATILITY AND CREATIVITY: THE MODEL IS CAPABLE OF GENERATING A WIDE RANGE OF IMAGES, FROM REALISTIC PHOTOGRAPHS TO ARTISTIC RENDERINGS, BASED ON THE SPECIFICITY AND CREATIVITY OF THE TEXT PROMPT.

AI AND MACHINE LEARNING: STABLE DIFFUSION REPRESENTS A SIGNIFICANT ADVANCEMENT IN AI AND MACHINE LEARNING, SHOWCASING THE POTENTIAL OF THESE TECHNOLOGIES IN CREATIVE FIELDS.

ACCESSIBILITY: THIS TECHNOLOGY IS BECOMING INCREASINGLY ACCESSIBLE, ENABLING MORE PEOPLE TO EXPERIMENT WITH AI-DRIVEN ART AND IMAGE GENERATION. I.E. DALI3 IN GPT4

STABLE DIFFUSION MODEL

WE ARE USING PRE-TRAINED MODEL FROM HUGGING FACE WEBSITE API NAMED AS STABLE-DIFFUSION-V2.

The screenshot shows the Hugging Face website interface. At the top, there is a navigation bar with links for Models, Datasets, Spaces, Docs, Solutions, Pricing, Log In, and Sign Up. The main content area displays the model card for 'stabilityai/stable-diffusion-2'. The card includes a 'Model card' tab, a 'Community' tab with 72 members, and sections for 'Edit model card', 'Downloads last month' (280,037), and an 'Inference API' section with a 'Text-to-Image' input field and a 'Compute' button. Below these are sections for 'Spaces using stabilityai/stable-diffusion-2' (503) and a list of contributing organizations. At the bottom, there are two generated images: one of a dog wearing a red beret and another of an astronaut riding a horse.

Hugging Face Search models, datasets, users...

Models Datasets Spaces Docs Solutions Pricing Log In Sign Up

s. stabilityai/stable-diffusion-2 like 1.68k

Text-to-Image Diffusers StableDiffusionPipeline stable-diffusion Inference Endpoints arxiv:2202.00512 arxiv:2112.10752 arxiv:1910.09700 License: openrail++

Model card Files and versions Community 72 Edit model card

Downloads last month
280,037

Inference API

Text-to-Image Your sentence here... Compute

This model can be loaded on the Inference API on-demand.

JSON Output Maximize

Spaces using stabilityai/stable-diffusion-2 503

- stabilityai/stable-diffusion
- guoyww/AnimateDiff
- darkstorm2150/protogen-web-ui
- society-ethics/DiffusionBiasExplorer
- multimodalart/civital-to-hf
- Rifd/ngees_doang

STABLE DIFFUSION MODEL

Step – 1 Setting up the environment

Step – 2 Importing Libraries

Step – 3 Combining Text and image in one result

Step-4 Configurations for Generation with Stable Diffusion

Step-5 Generating Text to Image.

STABLE DIFFUSION MODEL

1. Setting up the Environment

```
!pip install --upgrade diffusers transformers -q
```

```
[2] import torch  
print(torch.cuda.is_available())
```

False

```
[ ] !pip install torch torchvision torchaudio
```

```
!nvcc --version
```

```
nvcc: NVIDIA (R) Cuda compiler driver  
Copyright (c) 2005-2022 NVIDIA Corporation  
Built on Wed_Sep_21_10:33:58_PDT_2022  
Cuda compilation tools, release 11.8, V11.8.89  
Build cuda_11.8.r11.8/compiler.31833905_0
```

```
[ ] !pip install torch==1.12.0+cu113 torchvision==0.13.0+cu113 torchaudio==0.12.0+cu113 -f https://download.pytorch.org/wheel/torch\_stable.html
```

```
[ ] import torch  
print(torch.cuda.is_available())
```

True

```
[ ] print(torch.__version__)  
print(torch.version.cuda)
```

2.1.0+cu118
11.8

```
[ ] !pip install diffusers
```

STABLE DIFFUSION MODEL

3. Configurations for AI-Driven Image and Prompt Generation with Stable Diffusion and GPT-2

```
[ ] class CFG:  
    device = "cuda"  
    seed = 42  
    generator = torch.Generator(device).manual_seed(seed)  
    image_gen_steps = 50  
    image_gen_model_id = "stabilityai/stable-diffusion-2"  
    image_gen_size = (400,400)  
    image_gen_guidance_scale = 9  
    prompt_gen_model_id = "gpt2"  
    prompt_dataset_size = 31000  
    prompt_max_length = 12
```

Configuration of the Stable Diffusion.

STABLE DIFFUSION MODEL

The screenshot shows the Hugging Face Hub interface. On the left, there's a sidebar with a profile picture of a person with a yellow face, the text "Hugging Face", and a search bar "Search models, datasets, users...". The sidebar also lists "Profile", "Account", "Organizations", "Billing", "Access Tokens" (which is currently selected), and "SSH and GPG Keys". On the right, under "Access Tokens", there's a section titled "User Access Tokens" with a descriptive text about how they authenticate identity. Below this is a table showing a single token entry:

Text-2-Image Model	READ	Manage
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	Hide	Copy

At the bottom of the sidebar, there's a "New token" button.

4. Initializing Stable Diffusion Pipeline for Image Generation

```
[ ]  
image_gen_model = StableDiffusionPipeline.from_pretrained(  
    CFG.image_gen_model_id, torch_dtype=torch.float16,  
    revision="fp16", use_auth_token='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', guidance_scale=9  
)  
image_gen_model = image_gen_model.to(CFG.device)
```

```
vae/diffusion_pytorch_model.safetensors not found  
Keyword arguments {'guidance_scale': 9} are not expected by StableDiffusionPipeline and will be ignored.  
Cannot initialize model with low cpu memory usage because `accelerate` was not found in the environment. Defaulting to `low_cpu_mem_usage=False`. It is strongly recom  
```
```

```
pip install accelerate
```
```

```
.
```

```
Loading pipeline components... 100%
```

```
5/5 [00:39<00:00, 7.35s/it]
```

STABLE DIFFUSION MODEL

▼ 5.Image Generation Function Using Stable Diffusion Model

```
[ ] def generate_image(prompt, model):
    image = model(
        prompt, num_inference_steps=CFG.image_gen_steps,
        generator=CFG.generator,
        guidance_scale=CFG.image_gen_guidance_scale
    ).images[0]

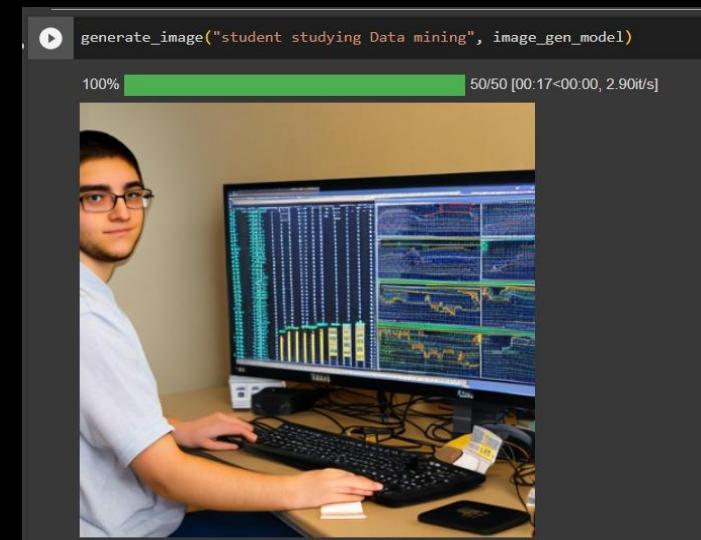
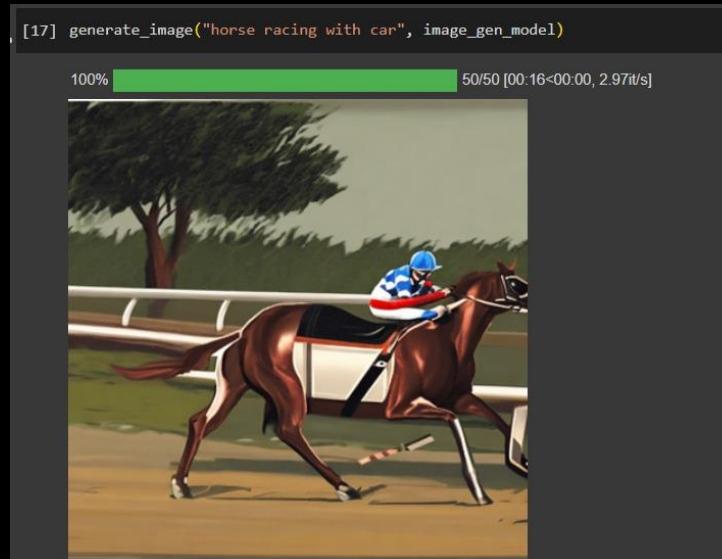
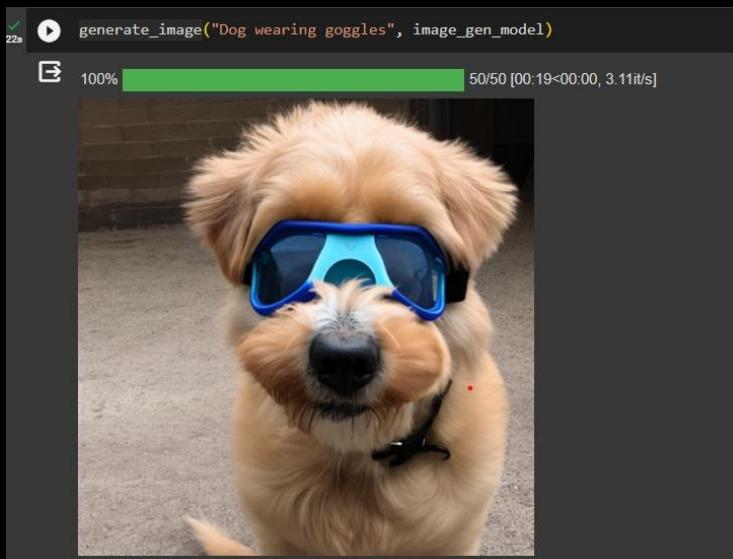
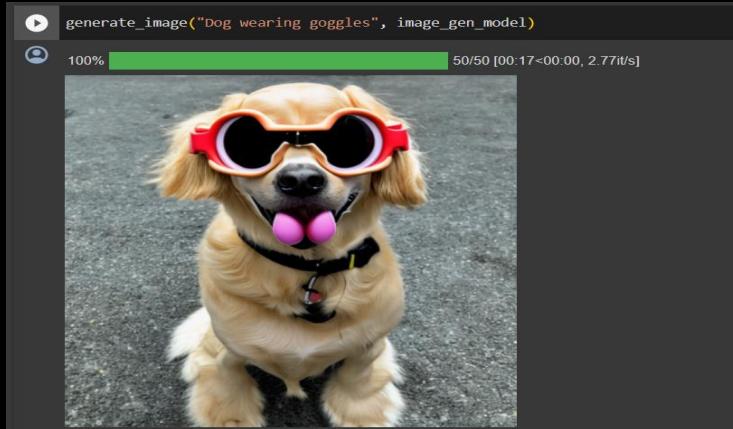
    image = image.resize(CFG.image_gen_size)
    return image
```

Outcome of our model



STABLE DIFFUSION MODEL

- Accuracy of our project is 75-80%
- out of 10 images 8 images came out to be similar what we have given input.
- In the 2nd run of the code the model is giving more accurate description in the input.
- As we train more data the model will give positive result.



PROBLEMS WE FACED

PROBLEMS WE FACED

```
[12] # Verify the content of captions
print(captions[:5]) # Print the first 5 captions to check their content

# Verify the type of each caption to ensure they are strings
print([type(caption) for caption in captions[:5]])

[nan, nan, nan, nan, nan]
[<class 'float'>, <class 'float'>, <class 'float'>, <class 'float'>]

▶ captions_df
[{"image_name": "image_name1", "comment_number": 1, "comment": "comment1"}, {"image_name": "1000092795.jpg", "comment_number": 0, "comment": "Two young guys with shaggy hair."}, {"image_name": "1000092795.jpg", "comment_number": 1, "comment": "Two young , White males are..."}, {"image_name": "1000092795.jpg", "comment_number": 2, "comment": "Two men in green shirts are..."}, {"image_name": "1000092795.jpg", "comment_number": 3, "comment": "A man in a blue shirt stand..."}, {"image_name": "...", "comment_number": "...", "comment": "..."}, {"image_name": "158911", "comment_number": 0, "comment": "A man in shorts and a Hawaii..."}, {"image_name": "158912", "comment_number": 1, "comment": "A young man hanging over the..."}, {"image_name": "158913", "comment_number": 2, "comment": "A man is leaning off of the ..."}, {"image_name": "158914", "comment_number": 3, "comment": "A man riding a small boat in..."}, {"image_name": "158915", "comment_number": 4, "comment": "A man on a moored blue and w..."}]
```

158916 rows × 3 columns

- NaN values removal
- Image not properly resized
- Model training .
- GAN Limitation
- GPU Problem
- Computation for processing images

REFERENCES

REFERENCES

1. High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ * Andreas Blattmann¹ * Dominik Lorenz¹ Patrick Esser Bjo rn Ommer¹

2. <https://github.com/Stability-AI/stablediffusion/blob/main/README.md>

3. <https://huggingface.co/stabilityai/stable-diffusion-2-1?text=dog+wearing+goggles>

4. <https://www.kaggle.com/datasets/adityajn105/flickr30k>

5.GAN lecture notes (A Brief Introduction to Generative Adversarial Networks (GAN))

DEMO

THANK YOU!

Under Guidance: Dr. Sunnie Chung

By Rahul Padwani
Jaykumar Mistry