

Harmonizing Fine-tuned Llama 2 for Content Generation with Stable Diffusion for Image Synthesis in Article Creation

Shenbagam P.¹, Thrisha Vaishnavi K S.², Hariprakassh S.³, Abhirami K.⁴, Abiram B.⁵, Rakesh Nandhaa K S.⁶

¹Assistant Professor, Department of Artificial Intelligence and data Science, Kumaraguru college of Technology, Coimbatore, India

^{2,3,4,5,6}UG students, Department of Artificial Intelligence and data Science, Kumaraguru college of Technology, Coimbatore, India

E-mail: ¹shenbagam.p.it@kct.ac.in, ²thrishavaishnavi@gmail.com, ³hariprakassh123@gmail.com, ⁴abhiramikathirvel@gmail.com, ⁵abiram.20ad@kct.ac.in, ⁶rakeshnandhaa.20ad@kct.ac.in

Abstract

The research explores the integration of generative AI in multimedia content production using a fine-tuned Llama 2 model for text generation and the Stable Diffusion algorithm for image synthesis. The research analyses the fine-tuned Llama 2-7b-chat model's adaptability to specific content generation contexts, enhanced by a unique dataset and QLoRa, a Quantized Low-Rank Adaptation for parameter-efficient fine-tuning, achieving significant reductions in training loss and nuanced quality in the generated content. Notably, the model's evaluation yielded an impressive perplexity score of 1.49, indicating advanced predictive performance. Additionally, stable diffusion's ability to transform textual descriptions into intricate images, highlighting its potential in AI-mediated content creation is demonstrated. The experiments and qualitative analyses reveal improvements in efficiency and creativity, emphasizing the collaborative potential of these models to revolutionize multidisciplinary content generation. The research underscores the transformative impact of fine-tuned generative models on content creation and offers insights into the broader implications for future AI research, while acknowledging the critical need for ethical considerations in the deployment of such technologies.

Keywords: Generative AI, Llama 2, Stable Diffusion, QLoRa, Parameter-Efficient fine-tuning

1. Introduction

In recent years, Generative Artificial Intelligence (AI) has emerged as a frontier technology, revolutionizing digital content creation and consumption. Central to this innovation are algorithms that learn from vast data sets to produce new, previously unseen outputs that mimic the original data in style and content. These models include generative adversarial networks (GANs), transformer models such as GPT variants, and diffusion models like Stable Diffusion, significantly expanding AI's creative capabilities [3]. This research focuses on the advancements in generative AI, particularly the fine-tuned Llama 2 model for text generation and the Stable Diffusion algorithm for image synthesis, highlighting significant enhancements in efficiency and creativity in AI-driven content production.

Significant strides have been made in the advancement of AI models for text generation and image synthesis. For example, OpenAI's GPT-3 has shown remarkable abilities in generating coherent and contextually relevant text across various applications, while models like DALL-E have demonstrated the ability to create detailed images from textual descriptions. However, the potential of combining the high-performing, fine-tuned Llama 2 model, renowned for its bespoke text generation, with the visually sophisticated Stable Diffusion, capable of rendering text descriptions into complex imagery, remains largely unexplored [1].

Within this investigation, the Llama 2 model, particularly its Llama 2-7b-chat variant, is scrutinized for its adeptness in adjusting to nuanced content creation parameters. The implementation of QLoRa, a Quantized Low-Rank Adaptation method, optimizes parameter efficiency, leading to significant reductions in training loss and improvements in the granularity and sophistication of the content produced [21]. For example, while prior studies such as on GPT-3 focused on general text generation capabilities, our approach fine-tunes the Llama 2 model specifically for article creation, yielding a training loss reduction to 0.66. Additionally, while models like DALL-E are used for generating images, our use of Stable Diffusion for seamless integration with text generation is a distinctive contribution. Case studies involving personalized marketing content and entertainment industry applications demonstrate the practical benefits and enhanced performance of our integrated approach.

The primary objective of this paper is to optimize the process of content generation using the fine-tuned Llama 2-7b-chat model. This involves fine-tuning the model on specific datasets and tasks relevant to article creation, enabling it to generate coherent and contextually relevant textual content with enhanced efficiency and accuracy. Another key objective is to implement Stable Diffusion techniques for image synthesis, allowing for the generation of high-quality images that align seamlessly with the generated textual content. This research provides a comprehensive examination of the synergies between text and image generative models, supported by robust scientific methodology, detailing dataset preparation, fine-tuning strategies, and iterative processes leading to a training loss of 0.66, demonstrating the effectiveness of our approach.

2. Literature Survey

The research showcases how language models, such as GPT-3, exhibit the ability to handle diverse tasks even with minimal task-specific data, showcases how fine-tuning can enhance performance on specific domains [11]. The study's findings provide a compelling argument for fine-tuning larger models like Llama 2 7b to refine their content generation capabilities.

The study from Training large scale neural language models dives deep into the training methodologies of neural language models, exploring the impact of large-scale datasets and architectural innovations. It emphasizes the importance of fine-tuning these models to enhance their generative capabilities and align them with specific tasks [12].

A key resource for understanding how latent diffusion models, akin to Stable Diffusion, can generate high-quality images is discussed in the research. The authors investigate the model's ability to synthesize images that accurately reflect given text descriptions, focusing on the balance between the fidelity of the generated images and the computational resources required [7].

The research provides deep insights into the generation of high-fidelity images. While focused on a different generative model, the paper's exploration of factors that influence image quality is relevant to understanding how Stable Diffusion synthesizes images from text descriptions [17]. The study offers methodologies that could potentially be applied to improve the image synthesis process and outcome in models like Stable Diffusion.

3. Content Generation using Fine Tuned Llama 2-7b-Chat-Model

3.1 Llama 2 Model Architecture

The Llama 2 model architecture represents a leap forward in the field of natural language processing. Built upon transformer-based technologies, it is designed to handle a vast array of language tasks. It is a decoder only model. Llama-2 is a set of three models rather than a single model. The quantity of parameters in each of these models is the only distinction between them. The Llama-2 models have parameters of 7B, 13B and 70B, in order of smallest to largest. There are many aspects that go into making LLaMA-2 effective [2].

While The Llama model underwent training with a context length of 2,000 tokens, The LLaMA-2 model is trained with an extended context length of 4,000 tokens. Furthermore, grouped query attention (GQA) is implemented by LLaMA-2 in every layer [9]. Based on an unlabelled textual corpus and the subsequent token prediction aim, all LLMs2 adhere to a (relatively) common and straightforward pre-training procedure. Given that the pre-training process for Large Language Models (LLMs) follows standardized procedures, the amount and quality of data used during pre-training significantly influence model performance. When pretraining an LLM with a larger and higher-quality dataset, especially for base models that have already undergone training, the resulting final model tends to exhibit improved performance. However, the models start with a pre-training process and a modified and enhanced model architecture. With an architecture designed for faster inference and pre-trained across a larger amount of data, LLaMA-2 enables the formation of a more comprehensive knowledge base. Additionally, Llama 2 is more effective due to the pretraining conditions provided below [15]. To handle the 2 Trillion tokens and internal weights, Llama2 uses a method known as Root Mean Square Normalisation, or RMSNorm. Llama 2 uses the SwiGLU activation function to decide whether a particular neuron should be active. Finally, Llama-2 uses a novel mathematical technique known as "Rotary Positional embedding," or RoPE, to make sure it comprehends the idea that the placements of words in sentences matter just as much as the words themselves.

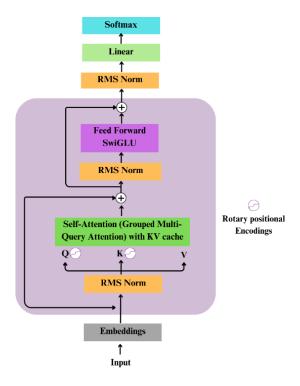


Figure 1. Optimized Architecture of Llama 2

The optimized architecture that is designed specifically for Llama 2 model along with the RMSNorm layer, SwiGLU activation function Layer and the integration of RoPE is represented in the above Figure 1 [14].

3.1.1 Llama 2-7B-Chat Model

Llama 2-7b-chat is a customized version of Llama 2 designed for dialog-based use cases. It is fine-tuned through Supervised Fine Tuning (SFT) and reinforcement learning from human feedback (RLHF). The pre-trained Llama 2 base model had been trained using supervised fine tuning to provide responses in the format that users are likely to expect in a chatbot [19]. Here the model is fine-tuned using prompt response pairs to minimize the difference between the model's own response for a given prompt and the example response given in the data provided. In RLHF, a "reward model" is trained with direct human feedback to identify patterns in the replies that people find most appealing. The reward model is then utilised to train Llama-2-chat through reinforcement learning, transforming its predictions into a scalar reward signal.

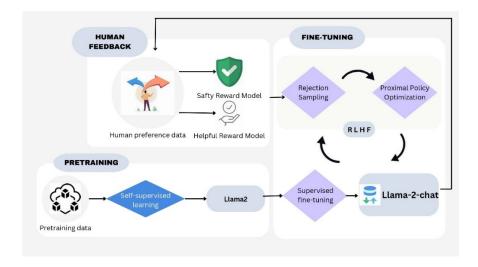


Figure 2. Llama 2 Chat Model Training

The procedure of training the Llama 2 conversation model making use of Supervised Fine Tuning and Reinforcement Learning through user responses is depicted in Figure 2 [6].

3.2 Fine Tuning of Llama 2-7B-Chat

3.2.1 Creating the Dataset

To fine-tune the Llama 2-7B-chat model, a specialized dataset of 150 rows was curated. This dataset consisted of high-quality dialogue pairs that target various domains of knowledge and conversational contexts. The data was collected from a wide array of articles available on the internet, supplemented by AI-generated content, covering a variety of subjects and topics. The choice of 150 rows ensures sufficient fine-tuning to enhance the model's contextual understanding and conversational adaptability while avoiding catastrophic forgetting, maintaining the balance between specificity and the model's foundational knowledge.

3.2.2 Parameter Optimized Fine Tuning

The ever-growing capabilities of Large Language Models (LLMs) like Llama 2 come at the cost of immense computational demands. Fine-tuning these models for specific tasks, like content generation for article creation, typically requires significant resources and time, hindering their practical application [5]. Parameter-Efficient Fine-tuning (PEFT) techniques aims to overcome these limitations by intelligently selecting and updating only a subset of LLM parameters relevant to the desired task. This approach leverages two key principles: Focus on Task-Specific Parameters, Reduced Precision for Non-Critical Parameters.

Among various PEFT methods, LoRA (Low-Rank Adaptor) and QLoRA (Quantized Low-Rank Adaptor) stand out for their effectiveness and unique approaches. QLoRa is Built upon LoRA, further enhances efficiency by quantizing the remaining parameters to lower precision levels. This reduces memory consumption and computational costs without sacrificing significant accuracy. LoRa replaces large weight matrices within specific model layers with a combination of smaller matrices and low-rank adaptors [16]. These adaptors capture task-specific information efficiently, resulting in significant parameter reduction.

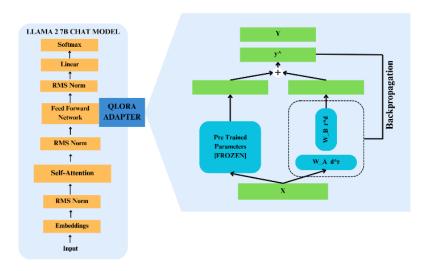


Figure 3. Architecture of Fine-tuned Llama 2-7B-Chat Model

The architecture of the fine-tuned Llama 2-7B Chat model using QLoRa and the decomposition of larger weight matrix into two smaller weight matrices along with the rank r and dimension of the model d in LoRa technique is represented in the above Figure 3. The rank parameter of the matrix controls the size of the smaller matrix.

3.2.3 Training and Loss

Fine-tuning the LLama2-7B-chat model with the QLoRa optimization technique involved an iterative process. Various parameter combinations were tested to identify the most effective configuration for enhancing model performance.

 Table 1. Parameters for Training

Parameters	Alpha	Rank	Dropout	Learning Rate
Values	16	64	0.1	2e-4

During the training phase, by using the parameters mentioned in the above Table 1, along with Adam optimizer, LLama2-7B-chat model underwent rigorous sequences of learning and evaluation to minimize this loss. To further enhance efficiency, the base model parameters were quantized to 4-bit precision using the "nf4" format [18]. This significantly reduces memory consumption and computational costs without using nested quantization. The model was trained over multiple epochs, where each epoch processed the entire dataset and updated the model's weights in response to the loss function. Throughout the training cycles, a continuous decrease in the loss metric was observed with a final loss of 0.24. This decline is reflective of the model's increasing capacity to comprehend and generate text that aligns closely with human conversation patterns. The fine-tuned model is saved in HuggingFace for further use.

3.2.4 Evaluation Metric

During the process of fine-tuning the Llama 2 model, we employed perplexity as the primary evaluation metric to measure the model's language understanding and predictive performance [20]. Perplexity, which gauges the model's uncertainty in predicting the next word in a sequence, serves as an inverse indicator of the model's performance. It is usually used only to determine how well a model has learned the training set. Lower values of perplexity indicate a more confident model with accurate predictions, whereas higher values like 10 to 100 suggest a model that is less certain about its predictions [8]. By optimizing our model to minimize perplexity, we ensured that the generated text was both contextually coherent and predictively precise, leading to more reliable content generation. Upon evaluating fine-tuning the Llama 2 model, We accomplished a perplexity score of approximately 1.49. This remarkably low perplexity indicates an exceptionally high level of predictive performance, affirming the model's ability to generate text with a high degree of certainty and fluency.

3.3 Process of Content Generation

To streamline the content creation pipeline, we incorporated LangChain, a tool designed to interface coherently with LLMs. Langchain helps LLM to access the external files and interact with external environments [13]. Leveraging LangChain, we devised a system where the user inputs are initially processed to understand the intent and context. The fine-tuned Llama 2 7b model, embracing its conversational prowess, then generates text that adheres to the specified themes or topics.

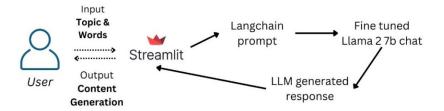


Figure 4. Process of Content Generation

Thus, by using LangChain, the output for content generation is created as shown in Figure 4.

3.4 Deployment using Streamlit

Streamlit's role in our architecture is pivotal for providing an interactive front end. The user interface created with Streamlit is intuitive, allowing for easy input of textual prompts and displaying the synthesized content. It enables real-time adjustments to the content generation parameters, offering users the ability to customize the output.

4. Image Generation using Stable Diffusion Model

4.1 Stable Diffusion Architecture

Stable Diffusion serves as a generative model specifically designed to create high-quality images based on textual descriptions. This technology harnesses the power of deep learning and a specific type of neural network known as a diffusion model. As opposed to traditional generative models, Stable Diffusion stands out for its ability to create detailed, coherent images by gradually refining noise into structured visuals over multiple iterative steps. It works by initially introducing random noise into an image canvas and then progressively denoising this input through the guidance of trained neural networks, which have learned to associate words with visual elements [4]. Stable Diffusion models are trained on large datasets of images and their associated textual descriptions, enabling them to generalize and produce new images that closely align with a wide range of descriptive prompts.

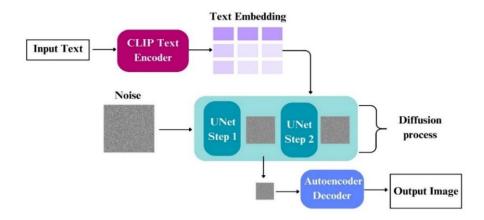


Figure 5. Architecture of Stable Diffusion

The architecture is split into three main components as described in the above Figure 5 [10] the variational autoencoder, which is used to decode the images from the latent space into the pixel space. The U-Net predicts denoised image representation of noisy latent. A transformer-based encoder called CLIP text encoder maps the input tokens to embeddings which is given as input to U-Net.

4.2 Process of Image Generation

Using a Stable Diffusion model typically involves the following steps:

Input Preparation: The process begins with a user providing a textual prompt. This prompt is passed through a CLIP text encoder to create an embedding that captures the semantic meaning of the text.

Noise Initialization: An initial noise image in latent space is generated, usually by sampling from a normal distribution. This image contains no discernible content and acts as a canvas onto which the model will project the user's textual description.

Diffusion Process: The noise image along with the text embedding is given as input to U-Net. It undergoes a series of transformations through the diffusion process. At each step, the model applies learned weights to predict and revert a small amount of the noise added. The embedding from the Transformer model is used to guide these transformations, ensuring that the refinements are aligned with the textual prompt.

Iterative Refinement: The refinement process is iterative, with the model gradually removing noise and adding structure to the image. Users can monitor intermediate stages to see how the model progresses from randomized noise to a coherent image.

Final Image Synthesis: Once the series of diffusion steps is complete, the result is a detailed image obtained through the decoder of variational autoencoder that visually interprets the textual prompt provided at the outset.

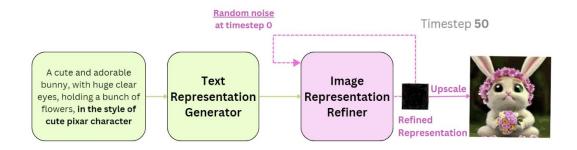


Figure 6. Example of image generation using Stable Diffusion

The above Figure 6 [22] showcases the gradual transformation from noise to final image through Stable Diffusion given a text prompt. The timestep can be adjusted as needed; increasing the number of timesteps generally results in a clearer image.

4.3 Deployment using Streamlit

To democratize access and facilitate interaction with the Stable Diffusion model, we utilize Streamlit, an open-source application framework tailored for machine learning and data science endeavors. Streamlit provides an effecient way for users to engage with the Stable Diffusion model, offering an intuitive interface where they can input textual prompts and receive generated images in real-time.

5. Integration of Content Generation and Image Generation

First, the topic for the content is provided as input to the fine-tuned LLama 2 7b chat model. The model then generates the desired content as output. Following this, the user manually provides a text prompt related to the generated content or the required content for the image generation. This prompt is used as input to the Stable Diffusion model to produce the final clear image related to the content.

6. Results and Discussion

The process fine-tuning of the Llama 2 7b chat model was meticulously carried out with a particular emphasis on reducing training loss and enhancing model response quality. Our results indicate a significant reduction in training loss to 0.24, demonstrating the effectiveness of the fine-tuning process. Utilizing a learning rate of 2e-4 and the Adam optimizer, the model exhibited improved conversational capabilities and increased contextual understanding. Most notably, the evaluation of the fine-tuned model yielded an impressive perplexity score, averaging around 1.49, which underscores the model's advanced predictive performance. This low perplexity is indicative of the model's proficiency in accurately predicting the next word in a sequence, reaffirming its enhanced generation capability. The combined results of diminished training loss and low perplexity point towards a highly optimized Llama 2 7b model that excels in generating precise, context-appropriate content tailored to our specific production needs.

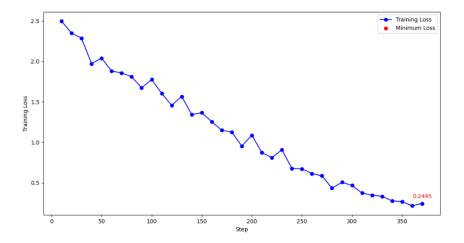


Figure 7. Training Loss Over Steps

A continuous decrease in the training loss is observed from the above Figure 7 over multiple steps and multiple epochs, thus resulting in an efficient fine-tuned Llama 2-7b-chat model.

Table 2. Comparative Analysis with Baseline Models

Model	Training Loss	Perplexity Score
Baseline Llama 2 7b	0.40	11.50
Baseline Llama 2 7b chat	0.35	10.30
Fine-tuned Llama 2 7b chat	0.24	1.49

The above Table 2 highlights the comparative analysis of the fine-tuned Llama 2 7b chat model with the baseline model. Stable Diffusion, on the other hand, served as the backbone for our image synthesis endeavors. We leveraged its state-of-the-art algorithm to convert text descriptions into high-fidelity images, complementing the textual content with visually appealing elements. The images produced were not only contextually relevant but also captured the intricate details conveyed by the text descriptions, thus demonstrating Stable Diffusion's potent capabilities.

7. Conclusion

In conclusion, our exploration has reinforced the value of fine-tuning generative models like Llama 2 for specialized content creation tasks. The training adaptations we incorporated have led to a model that performs with a substantial increase in contextual coherence, providing richer and more engaging conversational content generation. The effective combination of the fine-tuned Llama 2 7b chat model tailored through QLoRa with the visual synthesis prowess of Stable Diffusion holds the promise of revolutionizing the way multimedia content is produced, offering creators a powerful tool for crafting engaging and diverse narratives. This advancement was further complemented by LangChain's orchestration, which facilitated a seamless content generation workflow and ensured coherence and relevance in the output. The Streamlit interface provided an accessible and robust platform for real-time user interaction, bolstering the user experience and engagement.

Furthermore, the demonstrated efficiency and creative output establish a benchmark for future endeavors in AI-driven content creation. It also opens avenues for the personalization of generative AI to cater to a broad spectrum of creative requirements in various domains. While the current results are promising, we recognize that continued advancements and ethical

considerations are paramount to ensure responsible utilization and mitigate potential biases within AI-generated content.

Future studies could focus on further personalization of the models to cater to diverse creative domains, different languages extending their applicability and enhancing the finesse of the generated content. Moving forward, the aim is to further refine these technologies, expanding their capabilities and applications to foster new and innovative forms of digital storytelling and content creation. Additional research is also encouraged to address the ethical considerations of AI in content creation, ensuring responsible use and mitigating biases in AI-generated material.

References

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, et al." Attention Is All You Need". 31st International Conference on Neural Information Processing Systems (NeurIPS), no. 07 (2023): 6000–6010.
- [2] Touvron, Hugo, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov et al. "Llama 2: Open foundation and fine-tuned chat models." arXiv preprint arXiv:2307.09288 (2023).
- [3] Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, et al." Parameter-Efficient Fine-Tuning Methods for Pretrained Language Models: A Critical Review and Assessment". Nature Machine Intelligence, no. 05 (2023): 220-235.
- [4] Hu, Edward J., Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. "Lora: Low-rank adaptation of large language models." arXiv preprint arXiv:2106.09685 (2021).
- [5] Dettmers, Tim, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. "Qlora: Efficient finetuning of quantized llms." Advances in Neural Information Processing Systems 36 (2024).
- [6] Ho, Jonathan, Ajay Jain, and Pieter Abbeel. "Denoising diffusion probabilistic models." Advances in neural information processing systems 33 (2020): 6840-6851.

- [7] Robin Rombach, Andreas Blattmann, Dominik Lorenz, et al." High-Resolution Image Synthesis with Latent Diffusion Models".IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), no. 02 (2022): 10674-10685
- [8] Ning Ding, Yujia Qin, Guang Yang, et al." Parameter-efficient fine-tuning of large-scale pre-trained language models". Machine Intelligence, no. 05 (2023): 220-235.
- [9] Lermen, Simon, Charlie Rogers-Smith, and Jeffrey Ladish. "Lora fine-tuning efficiently undoes safety training in llama 2-chat 70b." arXiv preprint arXiv:2310.20624 (2023).
- [10] Pavlyshenko, Bohdan M. "Financial News Analytics Using Fine-Tuned Llama 2 GPT Model." arXiv preprint arXiv:2308.13032 (2023).1-14
- [11] Basile, Pierpaolo, Elio Musacchio, Marco Polignano, Lucia Siciliani, Giuseppe Fiameni, and Giovanni Semeraro. "LLaMAntino: LLaMA 2 models for effective text generation in Italian language." arXiv preprint arXiv:2312.09993 (2023).
- [12] Balachandran, Abhinand. "Tamil-Llama: A New Tamil Language Model Based on Llama 2." arXiv preprint arXiv:2311.05845 (2023).1-19
- [13] Pathak, Avik, Om Shree, Mallika Agarwal, Shek Diya Sarkar, and Anupam Tiwary. "Performance Analysis of LoRA Finetuning Llama-2." In 2023 7th International Conference on Electronics, Materials Engineering & Nano-Technology (IEMENTech), pp. 1-4. IEEE, 2023.
- [14] Bian, Junyi, Xiaolei Qin, Wuhe Zou, Mengzuo Huang, and Weidong Zhang. "Hellama: Llama-based table to text generation by highlighting the important evidence." arXiv preprint arXiv:2311.08896 (2023).
- [15] Xue, Zeyue, Guanglu Song, Qiushan Guo, Boxiao Liu, Zhuofan Zong, Yu Liu, and Ping Luo. "Raphael: Text-to-image generation via large mixture of diffusion paths." Advances in Neural Information Processing Systems 36 (2024).
- [16] Everaert, Martin Nicolas, Marco Bocchio, Sami Arpa, Sabine Süsstrunk, and Radhakrishna Achanta. "Diffusion in style." In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 2251-2261. 2023.

- [17] Tang, Raphael, Linqing Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus Stenetorp, Jimmy Lin, and Ferhan Ture. "What the daam: Interpreting stable diffusion using cross attention." arXiv preprint arXiv:2210.04885 (2022).
- [18] Zhan, Guanqi, Chuanxia Zheng, Weidi Xie, and Andrew Zisserman. "What Does Stable Diffusion Know about the 3D Scene?." arXiv preprint arXiv:2310.06836 (2023).
- [19] Stöckl, Andreas. "Evaluating a synthetic image dataset generated with stable diffusion." In International Congress on Information and Communication Technology, pp. 805-818. Singapore: Springer Nature Singapore, 2023. 805–818.
- [20] Sarafianos, Nikolaos, Xiang Xu, and Ioannis A. Kakadiaris. "Adversarial representation learning for text-to-image matching." In Proceedings of the IEEE/CVF international conference on computer vision, pp. 5814-5824. 2019.
- [21] Croitoru, Florinel-Alin, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. "Reverse Stable Diffusion: What prompt was used to generate this image?." arXiv preprint arXiv:2308.01472 (2023).
- [22] https://medium.com/polo-club-of-data-science/stable-diffusion-explained-for-everyone-77b53f4f1c4