# Dog Hip Angle Classification using Generated Images

**Sahil Kakkar**

Yeshiva University
skakkar@mail.yu.edu

**Abstract.** In this study, I aimed to enhance the accuracy of a classification model for dog hip X-ray images by incorporating generated images into the training dataset. The initial dataset comprised 2 classes: Big and Small, based on angles. To improve the performance, I fine-tuned a stable diffusion model to generate new X-ray images, which were then manually labeled and categorized based on their angles. Incorporating these labeled generated images into the training set and retraining the model resulted in a significant improvement. This demonstrates that generated images can effectively augment training datasets and enhance model performance. Extensive hyper-parameter tuning was performed to achieve these results, underscoring the potential of synthetic data in medical image classification tasks. Here's the GitHub link: https://github.com/Sahil1776/Dog-Hip-X-Ray-Classification-with-Generated-Images

## 1 Introduction

Accurate classification of medical images is crucial for diagnosing and treating various health conditions in both human and veterinary medicine. In veterinary orthopedics, assessing hip size and condition is essential for diagnosing issues such as hip dysplasia in dogs. Developing robust machine learning models for medical image classification often requires large, well-labeled datasets, which are challenging to obtain due to the specialized nature of the task and the time-intensive process of manual labeling. Recent advancements in generative models, such as stable diffusion, offer a promising solution to the data scarcity problem. These models can generate high-quality synthetic images that can augment existing datasets, potentially improving the performance of machine learning models. In this study, the application of stable diffusion-generated images was explored to enhance the classification accuracy of a model trained on dog hip X-ray images.

The primary dataset used in this study consisted of dog hip X-ray images categorized into two classes based on their size: Big and Small. Initial experiments with the classification model showed room for improvement. To address this, additional images were generated using a fine-tuned stable diffusion model. These images were then manually labeled and incorporated into the training dataset. Retraining the model with this augmented dataset led to a notable increase in accuracy. This study demonstrates the potential of synthetic data to enhance the performance of medical image classification models. By leveraging generated images, it is possible to improve model accuracy and reduce reliance on large, manually labeled datasets. The findings underscore the value of data augmentation techniques in medical imaging and present a new approach to overcoming the challenges associated with limited data availability.

## 2 Related Work

The project delves deeply into the augmentation of a dataset specifically designed for the classification of canine hips, harnessing the innovative prowess of advanced text-to-image generation models. This ambitious and forward-thinking undertaking capitalizes on the transformative capabilities of state-of-the-art generative models, which have the remarkable ability to fabricate highly detailed and intricate hip X-Ray images. These images portray dogs' hips, and their respective angles. Here are the main methods that are related to tasks performed in our study:-

### 2.1 Generative Adversarial Networks (GANs)

Generative Adversarial Networks (GANs) are highly relevant and beneficial in the context of the information provided above. Here's how GANs are helpful:

1. **High-Quality Image Synthesis**: GANs are known for their ability to generate high-quality and realistic images. In the context of our project, which aims to augment a dataset for the classification of canine hip conditions, GANs can be instrumental in synthesizing detailed and accurate X-ray images of dogs' hips. This is particularly valuable as the quality of synthetic images directly impacts the effectiveness of the dataset augmentation process.

2. **Diverse Image Generation**: One of the strengths of GANs is their ability to produce a wide variety of images. This diversity is crucial for our task, as we need to generate images that represent the two distinct categories of canine hips (big and small). By leveraging GANs, we can ensure that our augmented dataset encompasses a broad range of scenarios, thereby improving the robustness and generalization capabilities of our classification models.

3. **Efficient Hyperparameter Tuning**: GANs, like other generative models, require meticulous hyperparameter tuning to achieve optimal performance. The process of fine-tuning GANs involves adjusting parameters such as the learning rate, batch size, and network architecture to enhance image synthesis quality. This iterative experimentation and empirical validation process aligns perfectly with our methodology, as we strive to identify the best hyperparameter settings for our generative models.

4. **Leveraging Advanced Computational Infrastructure**: Training GANs is computationally intensive and demands robust hardware and software infrastructure. Our project benefits from state-of-the-art

computational resources, including high-performance GPUs and parallel processing capabilities, which are essential for efficiently training GANs. This infrastructure supports the seamless execution of the computationally demanding tasks associated with GAN training and image synthesis, accelerating the overall research progress.

## 2.2 Progressive and Style-Based GANs

Progressive and Style-Based GANs are highly relevant and beneficial in the context of the information provided above, particularly for the augmentation of a dataset tailored for the classification of canine hip conditions. Here's how these specific types of GANs are helpful:

a) **Progressive GANs**

1. **Gradual Improvement in Image Quality**: Progressive GANs (PGANs) are designed to improve image quality gradually during the training process by starting with low-resolution images and progressively increasing the resolution. This approach helps in synthesizing highly detailed and accurate X-ray images of dogs' hips. In the context of our project, this gradual improvement ensures that the generated images become more realistic over time, which is crucial for effective dataset augmentation.

2. **Stable Training**: Training GANs can be challenging due to issues like mode collapse and instability. Progressive GANs address these issues by focusing on lower resolutions first and then progressively increasing the complexity. This leads to more stable training and higher-quality image generation, which is beneficial for creating realistic and diverse images of canine hips in big and small categories.

3. **High-Fidelity Image Synthesis**: The ability of Progressive GANs to produce high-fidelity images is essential for our project. High-quality images ensure that the augmented dataset accurately represents the variations seen in real-world scenarios, thereby enhancing the robustness and generalization capabilities of our classification models for canine hip conditions.

b) **Style-Based GANs (StyleGANs)**

1. **Fine-Grained Control over Image Features**: StyleGANs introduce the concept of style vectors at different layers of the generator network, allowing for fine-grained control over various image features. In our project, this capability is particularly useful for synthesizing detailed and specific features of canine hips, such as bone structure and joint characteristics. This level of control can lead to more accurate and diverse synthetic images.

2. **Generation of High-Resolution Images**: StyleGANs are known for their ability to generate high-resolution images with remarkable detail. This is critical for our dataset augmentation task, as high-resolution X-ray images of dogs' hips are necessary for precise classification into big and small categories. The improved resolution contributes to better training data for our models, leading to more accurate diagnostic tools.

3. **Enhanced Variability and Realism**: By manipulating style vectors, StyleGANs can produce images with a wide range of variations, which is essential for creating a comprehensive and varied dataset. This enhanced variability helps in training models that can generalize better to new and unseen data, thereby improving the reliability and effectiveness of the classification models for canine hip conditions.

4. **Efficient Hyperparameter Tuning**: Both Progressive and Style-Based GANs require careful hyperparameter tuning to achieve optimal performance. The structured approach of these models, with their progressive layers and style vectors, allows for more systematic and efficient hyperparameter exploration. This aligns with our methodology of meticulous experimentation and empirical validation to identify the best settings for high-quality image synthesis.

## 2.3 Diffusion and Autoencoder-based Frameworks

Diffusion and Autoencoder-based frameworks play a crucial role in augmenting our dataset for the classification of canine hip conditions, specifically for generating high-quality and diverse X-ray images of dogs' hips categorized as big or small. Here's how these frameworks are helpful:

a) **Diffusion-Based Frameworks**

1. **Gradual and Detailed Image Generation**: Diffusion-based frameworks, such as Denoising Diffusion Probabilistic Models (DDPM), generate images through a process of gradually adding and then removing noise. This step-by-step approach allows for the creation of highly detailed and realistic images. In the context of our project, diffusion models can effectively generate intricate X-ray images of canine hips, ensuring that the synthetic images are of high quality and closely resemble real-world data.

2. **Robustness to Noise**: Diffusion models are inherently designed to handle noise, which makes them robust and stable during training. This robustness is beneficial for generating clear and accurate images of canine hips, even when dealing with complex structures and varying conditions. The stability during training reduces the risk of artifacts or unrealistic features in the generated images, enhancing the overall quality of the dataset.

3. **Enhanced Diversity and Realism**: The iterative nature of diffusion models allows for the generation of a wide variety of images by controlling the diffusion process parameters. This ability to produce diverse images is crucial for our dataset augmentation, as it ensures that the synthetic images cover a broad spectrum of hip conditions, thereby improving the generalization capabilities of our classification models.

b) **Autoencoder-Based Frameworks**

1. **Dimensionality Reduction and Reconstruction**: Autoencoder-based frameworks are designed to learn efficient representations of data by compressing it into a lower-dimensional latent space and then reconstructing it back to the original space. This capability is useful for generating synthetic X-ray images of canine hips by learning the underlying features and patterns in the data. Autoencoders can effectively capture the essential characteristics of big and small hip conditions, enabling the generation of realistic images.

2. **Noise Reduction and Denoising**: Autoencoders, particularly Denoising Autoencoders (DAEs), are adept at removing noise from input images. This feature can be leveraged to enhance the quality of synthetic images by reducing artifacts and improving clarity. In our project, denoising autoencoders can help produce clean and accurate X-ray images of canine hips, which are essential for reliable classification.

3. **Anomaly Detection**: Autoencoders can also be used for anomaly detection by comparing the reconstructed images with the

original ones. This feature is valuable for identifying and correcting any unrealistic or erroneous images generated during the augmentation process. Ensuring that only high-quality synthetic images are added to the dataset helps maintain the integrity and effectiveness of the classification models.

c) **Leveraging Advanced Computational Infrastructure**

1. **Handling Complex Computations**: Both diffusion and autoencoder-based frameworks require significant computational resources for training and image generation. Our advanced computational infrastructure, equipped with high-performance GPUs and parallel processing capabilities, is well-suited for handling these complex computations. This infrastructure allows us to efficiently train models and generate high-quality synthetic images at scale, accelerating the research process.

2. **Scalability and Efficiency**: The scalability of our computational resources ensures that we can easily expand our capabilities to accommodate the demands of training diffusion and autoencoder-based models. Whether it involves processing larger datasets or experimenting with more complex model architectures, our infrastructure is designed to support these needs seamlessly, ensuring efficient utilization of resources.

d) **Enhancing Dataset Augmentation**

1. **Improving Training Data Quality**: By leveraging diffusion and autoencoder-based frameworks, we can generate high-quality synthetic images that enrich our training dataset. The enhanced quality and diversity of these images contribute to more effective model training, resulting in better performance and accuracy in classifying canine hip conditions.

2. **Increasing Dataset Size and Variability**: The ability to generate a large number of realistic and varied images using these frameworks significantly increases the size and variability of the training dataset. This increase in data helps improve the robustness and generalization capabilities of the classification models, ensuring they perform well on real-world data.

In summary, diffusion and autoencoder-based frameworks are highly beneficial in our project for generating high-quality, diverse, and realistic X-ray images that augment our dataset for the classification of canine hip conditions. Their capabilities in detailed image generation, noise reduction, and anomaly detection, combined with our advanced computational infrastructure, make them essential components of our research methodology.

### 2.3.1 *Stable Diffusion Dreambooth*

Stable Diffusion DreamBooth is an advanced framework used for generating high-quality and contextually relevant images through fine-tuning of diffusion models. Here's how it can be specifically helpful for augmenting a dataset of X-ray images of canine hips, categorized as big or small:

1. **High-Quality Image Synthesis**: Stable Diffusion DreamBooth leverages the diffusion model's ability to produce detailed and realistic images. By fine-tuning this model on our dataset of canine hip X-rays, we can generate high-quality synthetic images that closely resemble real-world X-rays. This high-quality synthesis is crucial for ensuring that the augmented dataset accurately reflects the characteristics of canine hips.

2. **Contextual Adaptation**: DreamBooth allows for the adaptation of the diffusion model to specific contexts or domains. In our case, it means tailoring the model to understand and generate X-ray images of canine hips with different conditions (big or small). This contextual adaptation helps in producing images that are highly relevant and useful for our classification task.

In summary, Stable Diffusion DreamBooth provides a powerful framework for augmenting our dataset of canine hip X-rays. Its capabilities in high-quality image synthesis, contextual adaptation, enhanced diversity, and efficient fine-tuning make it an invaluable tool for improving the dataset and training robust classification models.

## 3  Methods

### 3.1  *Dataset*

We initially started with a dataset of images that were each accompanied by their specific captions. These captions were not just any captions; they were carefully crafted by experts in the field. The expertise involved ensured that the quality of these captions was extremely high. Each caption was meticulously written to accurately describe the corresponding image, making the dataset valuable for our purposes.

However, despite the high quality of the dataset, we encountered a significant limitation: the quantity of data. The dataset was simply not large enough to fulfill our needs, particularly because we required a diverse set of X-rays of dog hips. The challenge was that X-rays of dog hips taken from different angles are relatively rare. We needed a variety of angles to ensure comprehensive coverage, but the existing dataset did not provide this.

To overcome this limitation, we turned to advanced technology. We utilized the stable diffusion model, a powerful tool for generating images. This model allowed us to create new images that we could add to our existing dataset. By generating these new images, we were able to address the gap in our dataset, particularly by providing the missing X-rays with different angles. The stable diffusion model played a crucial role in expanding our dataset to better meet our needs.

As we generated these new images, we understood the importance of keeping our dataset well-organized. To achieve this, we created a metadata.csv file. This file served as a catalog for our expanded dataset. It contained the updated captions that corresponded to the newly generated images. Additionally, the metadata.csv file included the names of the image files, ensuring that each caption was properly linked to its respective image.

The creation of the metadata.csv file was a crucial step in maintaining the integrity and usability of our dataset. It allowed us to seamlessly integrate the newly generated images with the existing ones, ensuring that everything was well-organized and easy to access. By doing so, we were able to enhance our dataset, making it more comprehensive and better suited to our requirements. This systematic approach to expanding our dataset was essential in achieving our goals.

### 3.2  *Image Generation*

To augment the dataset, we embarked on the process of fine-tuning the stable diffusion model using our initial dataset. Fine-tuning is a

crucial step in machine learning, where the model's parameters are adjusted to better understand the specific features and structures that are unique to the data it's being trained on. In our case, this involved training the model to recognize the distinct characteristics of dog hip X-ray images. By refining the model's parameters, we were able to make the model more attuned to the specific patterns, shapes, and structures that are found in these X-rays.

After fine-tuning, the stable diffusion model was employed to generate a large number of new images. The generation process was meticulous, and the images produced were subjected to visual inspection. This inspection was necessary to ensure that the newly generated images maintained a high level of quality and were consistent with the original dataset. We wanted to make sure that the generated images not only resembled the original X-rays but also captured the variations in angles and structures that are critical for our analysis.

Through this process, we generated a substantial number of images—100,000 in total. This massive augmentation of our dataset provided us with a much richer and more diverse set of images to work with. However, generating images was only one part of the process. To further enhance the value of the dataset, many of these generated images were manually labeled. Manual labeling is a labor-intensive process but crucial for ensuring that the dataset is accurate and reliable. Each image was carefully reviewed, and labels were assigned to accurately describe the features and angles depicted.

In addition to manual labeling, we also utilized our previous model to produce predictions on the newly generated images. This involved running the generated images through the model to see how well it could predict the relevant features. However, predictions made by the model were not always perfect, so corrections were made by hand. This correction process was critical in refining the dataset, as it allowed us to combine the strengths of machine learning with human expertise. By correcting the predictions, we were able to further augment our dataset, making it even more robust and useful for training purposes.

This augmented dataset was then used to improve our previous model. The goal was to leverage the expanded dataset to increase the model's accuracy. With more data to learn from, the model could better generalize and make more accurate predictions when applied to new X-ray images. The diverse set of images, particularly with different angles, played a significant role in enhancing the model's performance.

Each image in our dataset had a corresponding caption that provided details about the image. These captions were not generic; they were specifically crafted to describe various aspects of the image, including the angles at which the X-rays were taken. This level of detail in the captions was essential for guiding the generation of images with varying angles, ensuring that our dataset covered a wide range of perspectives.

Some of the generated images, which illustrate the success of our augmentation process, are shown in Figure 1. These images demonstrate the variety and quality of the generated dataset, highlighting the effectiveness of our fine-tuning and generation efforts. By combining advanced modeling techniques with careful inspection and manual correction, we significantly enhanced our dataset, ultimately contributing to improved model accuracy.
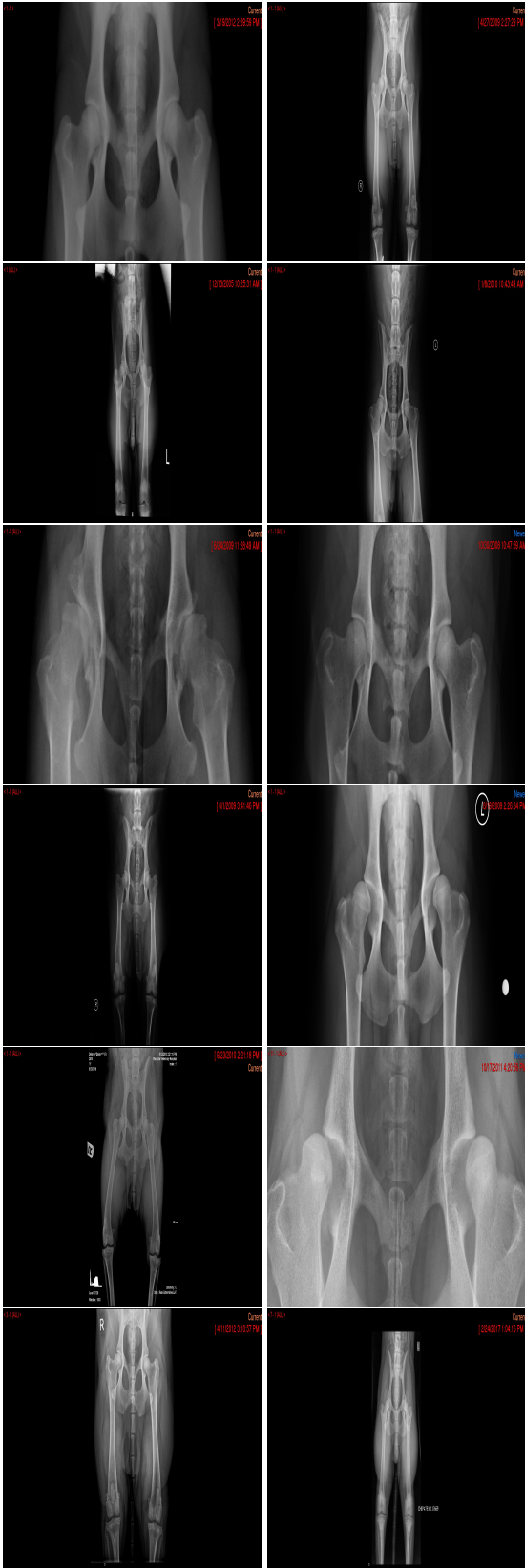


**Figure 1**: Generated Images of Dog's hips with varying angles

## 3.3  Model Architecture

The Stable Diffusion model is a prominent generative model designed to create high-quality images from latent space representations. Its architecture builds on principles from diffusion models and utilizes innovative techniques to achieve state-of-the-art image synthesis. Here's a detailed overview of its architecture:

1. **Diffusion Model Foundation**
Stable Diffusion is built upon the foundational principles of diffusion models, which are a class of generative models designed to create structured data from random noise. The core idea behind diffusion models is to start with a simple, unstructured input, typically Gaussian noise, and then gradually transform it into complex and meaningful data. This transformation is achieved through a step-by-step process that involves both the degradation and reconstruction of the data. The entire process can be broken down into two key components: the forward process, known as diffusion, and the reverse process, known as denoising.

The first component, the forward process or diffusion, is where the data undergoes a controlled degradation. This process begins with the original data, which could be an image, a sound, or any other form of structured information. The diffusion process then progressively adds Gaussian noise to the data over a series of steps. With each step, the data becomes more corrupted, losing its original structure and detail, until it eventually resembles pure noise. The purpose of this process is to map the data into a noisy space where it is indistinguishable from random noise. By the end of the diffusion process, the data has been completely transformed into a form of noise, erasing all recognizable patterns and structures.

The second component is the reverse process, or denoising, which is where the real magic of diffusion models happens. In this phase, a neural network is trained to reverse the effects of the diffusion process. The network learns to take the noisy data generated during the diffusion phase and gradually recover the original structured data. This denoising process is also performed in a series of steps, mirroring the steps of the forward process but in reverse order. The neural network begins with the noisy data and, step by step, removes the noise while reconstructing the original data. The goal of this process is to train the network to accurately predict and reverse the noise at each stage, ultimately restoring the data to its original, uncorrupted state.

Through the combination of these two processes, diffusion and denoising, Stable Diffusion is able to generate high-quality data from simple Gaussian noise. The forward process ensures that the data is thoroughly randomized, while the reverse process carefully reconstructs the original data from this noise. The success of diffusion models like Stable Diffusion lies in their ability to learn complex patterns and structures, allowing them to generate realistic and detailed outputs even from an entirely random starting point. This approach has proven to be highly effective in various applications, including image generation, where Stable Diffusion can create detailed and coherent images from noise.

2. **Latent Space Representation**
Unlike some diffusion models that directly manipulate data in the image space, Stable Diffusion takes a different approach by operating in a latent space. This method involves working with a compressed representation of the original data rather than the data itself, leading to increased efficiency and scalability. The process of operating in latent space can be broken down into several key steps, each contributing to the overall effectiveness of the model.

The first critical component in this approach is the Encoder Network. The encoder's role is to take the high-dimensional images, which consist of a large amount of data, and compress them into a lower-dimensional latent space representation. This compression process involves capturing the essential features and structures of the images while discarding extraneous details that are not crucial for the task at hand. By reducing the dimensionality of the data, the encoder significantly decreases the computational complexity of the subsequent steps. Instead of working with the full, high-dimensional image data, the model only needs to process these compact latent representations, making the entire process more manageable and efficient.

Once the data is compressed into latent space, the next step is the Latent Space Diffusion process. In this stage, the principles of diffusion, as described earlier, are applied not to the original high-dimensional images but to the latent representations produced by the encoder. This is a crucial distinction that sets Stable Diffusion apart from other models that operate directly on image data. By working in the latent space, the model can achieve the same transformative effects—gradually adding and then removing noise to generate new data—while using far fewer computational resources. This approach allows Stable Diffusion to handle larger datasets and more complex tasks with greater efficiency.

The diffusion process in latent space follows the same two key phases as in image space: the forward process (diffusion) and the reverse process (denoising). During the forward process, Gaussian noise is added to the latent representations step by step, gradually corrupting the structured data until it reaches a state of pure noise. Then, in the reverse process, a neural network is trained to reconstruct the original latent representations from the noisy ones, effectively reversing the corruption and recovering the structured data.

The benefits of operating in latent space are significant. Because the data has been compressed, the model can perform the diffusion process much more quickly and with less computational power. Additionally, the reduced dimensionality of the latent space allows the model to focus on the most critical features of the data, improving the efficiency and scalability of the diffusion process. This makes Stable Diffusion particularly well-suited for applications where computational resources are limited, or where large-scale data generation is required.

In summary, by using an encoder network to compress images into a latent space and then applying the diffusion process within this space, Stable Diffusion achieves a balance between performance and efficiency. This innovative approach allows the model to generate high-quality data while operating in a more scalable and resource-efficient manner, making it a powerful tool for various generative tasks.

3. **U-Net Architecture**
The core component of Stable Diffusion's architecture is the U-Net, a specialized type of neural network that is particularly well-suited for image-to-image tasks. U-Net has become a popular architecture in the field of deep learning due to its ability to handle complex image processing tasks with precision and efficiency. In the context of Stable Diffusion, the U-Net plays a crucial role in transforming latent representations into high-quality images. The architecture of the U-Net can be broken down into three main components: the

encoder path, the decoder path, and the skip connections, each serving a specific purpose in the overall process.

The first component is the Encoder Path. The encoder path is responsible for capturing high-level features from the input latent representations. This process involves passing the latent data through a series of convolutional layers, each followed by a downsampling operation. The convolutional layers help the network to extract important features, such as edges, textures, and patterns, from the input data. As the data progresses through the encoder path, it is gradually compressed, and its dimensionality is reduced. This allows the network to focus on the most critical high-level features while discarding less important details. The downsampling operations further reduce the resolution of the data, enabling the network to capture broader context and more abstract features at each stage.

The second component is the Decoder Path, which is tasked with reconstructing the latent space representations into detailed output images. In contrast to the encoder path, the decoder path gradually upsamples the data, increasing its resolution and adding finer details as it moves through the network. The decoder uses a series of upsampling operations to achieve this, effectively reversing the compression that occurred in the encoder path. At each stage of the decoder, the network refines the data, reconstructing the complex structures and details that make up the final image. The goal of the decoder path is to produce an output that is both accurate and detailed, closely resembling the original high-dimensional data from which the latent representations were derived.

A key feature of the U-Net architecture is its Skip Connections. These connections link corresponding layers in the encoder and decoder paths, allowing information to flow directly between them. Skip connections play a crucial role in preserving spatial information that might otherwise be lost during the downsampling and upsampling processes. By allowing the network to combine high-level features captured by the encoder with detailed local information retained in the skip connections, the U-Net is able to produce more accurate and realistic images. The skip connections ensure that the decoder has access to both the abstract, high-level features and the fine-grained details necessary for producing high-quality outputs. This combination is essential for generating images that are not only visually coherent but also rich in detail.

In the context of Stable Diffusion, the U-Net architecture is an integral part of the model, enabling it to generate high-quality images from latent space representations efficiently. The encoder and decoder paths work together to process and reconstruct the data, while the skip connections ensure that important spatial information is preserved throughout the process. This results in images that are both detailed and consistent with the input data.

A simplistic representation of the U-Net model architecture, which forms a key part of the Stable Diffusion model, is shown in Figure 2. This figure illustrates how the encoder path compresses the data, the decoder path reconstructs it, and the skip connections facilitate the transfer of critical information between the two paths. By leveraging the strengths of the U-Net architecture, Stable Diffusion is able to achieve impressive results in generating detailed and high-quality images from latent space.

## 4. Conditioning Mechanisms

Stable Diffusion incorporates advanced conditioning mechanisms to guide the image generation process, allowing for more control and precision in the outputs it produces. These conditioning mechanisms enable the model to generate images that are not only visually coherent but also aligned with specific requirements or descriptions provided by the user. There are two primary types of conditioning mechanisms integrated into Stable Diffusion: text-to-image conditioning and conditional inputs. Each of these mechanisms plays a vital role in shaping the final generated images, ensuring they meet the desired criteria.

The first type of conditioning mechanism is Text-to-Image Conditioning. This approach allows the model to generate images based on textual descriptions, effectively translating words into visual representations. The process begins with a text encoder, such as CLIP (Contrastive Language-Image Pretraining), which is designed to process and understand text input. The text encoder takes the provided textual description and converts it into embeddings—numerical representations that capture the semantic meaning of the text. These embeddings serve as a guide for the diffusion model, influencing the generation process at each step. By using text-to-image conditioning, Stable Diffusion can produce images that closely match the concepts and details described in the text. This capability is particularly useful in applications where users need to generate specific images based on verbal or written instructions, such as in creative design, content creation, or customized image generation.

The second conditioning mechanism involves Conditional Inputs. In addition to textual descriptions, Stable Diffusion can incorporate other types of conditions to further steer the image synthesis process. These conditions can include class labels, specific attributes, or any other predefined criteria that the user wants the generated image to adhere to. For example, a user might provide a class label that indicates the type of object or scene they want to generate, or they might specify certain attributes like color, style, or composition. The model integrates these conditional inputs into the diffusion process, ensuring that the generated images align with the specified conditions. This flexibility allows Stable Diffusion to produce highly targeted and customized images, making it suitable for a wide range of applications where precise control over the output is required.

By combining text-to-image conditioning with additional conditional inputs, Stable Diffusion offers a powerful and versatile image generation framework. The text-to-image conditioning provides a way to generate images based on rich and nuanced textual descriptions, while the conditional inputs allow for further refinement and customization of the generated images. Together, these conditioning mechanisms enable users to achieve highly specific and detailed outcomes, making Stable Diffusion an invaluable tool for tasks that demand both creativity and precision.

Through these conditioning mechanisms, Stable Diffusion empowers users to guide the image generation process with remarkable accuracy. Whether it's generating an image that perfectly matches a written description or creating an image that meets specific criteria, the model's ability to incorporate multiple forms of conditioning ensures that the final outputs are not only visually appealing but also aligned with the user's intentions. This integration of conditioning mechanisms is a key feature that sets Stable Diffusion apart, enabling it to deliver high-quality, customized images that cater to a wide variety of needs and applications.

## 5. Training Objectives

The training of Stable Diffusion involves optimizing a loss function that balances the reconstruction quality and adherence to the conditioning inputs:

a) Reconstruction Loss: This component of the loss function ensures that the model can accurately reconstruct the original data from its noisy latent representations. The goal is to minimize the difference between the generated images and the original images, thereby ensuring that the model preserves the essential features and details during the image generation process.

b) Conditioning Loss: This part of the loss function ensures that the generated images align with the specified conditions, such as textual descriptions or class labels. By incorporating conditioning loss, the model is guided to produce images that not only resemble the original data but also meet the requirements or attributes specified by the conditioning inputs. This helps in generating images that are consistent with the desired outputs based on the provided conditions.

### 6. Scalability and Efficiency

Stable Diffusion is designed to be scalable and efficient:

a) Efficient Training: One of the primary advantages of Stable Diffusion is its operation in latent space, which significantly reduces the computational burden compared to working directly with high-resolution images. By processing data in a compressed latent representation, the model can perform training and image generation more quickly and with less computational resource consumption. This approach allows for faster training times and more efficient generation of images.

b) High-Resolution Outputs: Despite operating in the latent space, Stable Diffusion is capable of generating high-resolution images. This is achieved by progressively refining the latent representations through the diffusion process. The model's architecture enables it to produce detailed and high-quality images, even though the initial computations are performed in a lower-dimensional space.

### 7. Inpainting and Editing Capabilities

Stable Diffusion models often include features for inpainting and editing images:

a) Inpainting: The inpainting feature allows users to modify specific parts of an image while preserving the rest of the image intact. This capability is useful for correcting artifacts, filling in missing parts, or adding new details to images. By focusing on particular regions of an image, inpainting helps to improve the overall quality and relevance of the generated content.

b) Editing: The editing capabilities of Stable Diffusion provide tools to make changes to existing images based on new conditions or prompts. This allows users to adapt and refine images according to updated requirements or creative directions, offering enhanced flexibility in image manipulation and generation.

### Leveraging Computational Infrastructure

1. **Parallel Processing**: Training Stable Diffusion models involves managing large volumes of data and complex computations. High-performance GPUs and parallel processing capabilities are essential for efficient training and image generation. These technologies enable the model to handle extensive data sets and perform computations concurrently, accelerating the training process.

2. **Handling Large Models**: Stable Diffusion models can be large and complex, demanding substantial memory and processing power. Efficient handling of these large models requires robust computational infrastructure that supports smooth training and generation processes. Adequate resources ensure that the model can perform ef-

fectively without bottlenecks or performance issues.

### Enhancing Dataset Augmentation

1. **Generating High-Quality Images**: The model's ability to create realistic and high-resolution images makes it an excellent tool for augmenting datasets with synthetic examples. These high-quality images can enhance the diversity and richness of datasets, contributing to more robust training and improved model performance.

2. **Customizing Outputs**: The conditioning mechanisms of Stable Diffusion enable the generation of images that meet specific requirements or attributes. This feature is valuable for creating diverse and tailored datasets, as it allows for precise control over the characteristics of the generated images, aligning them with particular needs or specifications.

In summary, the Stable Diffusion model combines advanced diffusion principles with efficient latent space representations and powerful conditioning mechanisms. Its architecture enables high-quality image synthesis, scalability, and flexibility, making it a valuable tool for augmenting datasets and enhancing various image generation tasks.
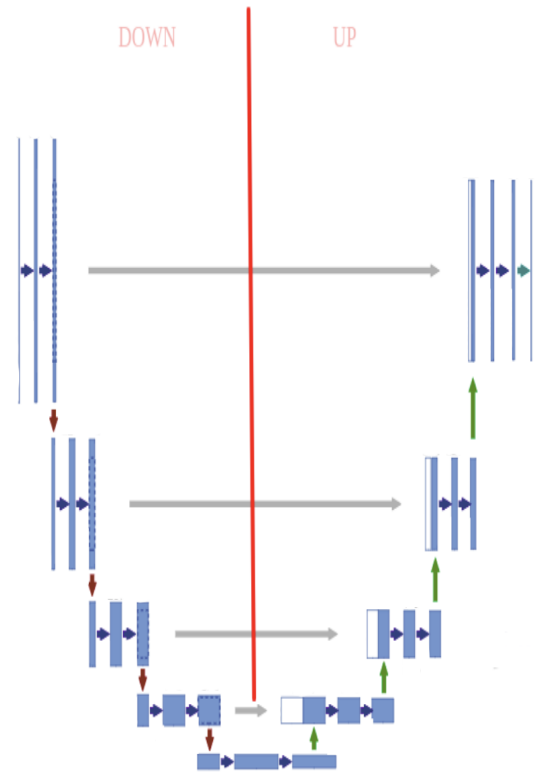


**Figure 2**: U-Net pipeline

## 4 Results

Initially, the performance of the model in terms of correlation and accuracy was not satisfactory. The initial results highlighted that the model struggled to effectively learn from the existing dataset, leading to suboptimal predictions and a need for improvement. To address this challenge, a substantial number of synthetic images were generated—specifically, 10,000 images were created to serve as additional data for model training.

These generated images were utilized in several ways to enhance the model's performance. Predictions were made on these synthetic images, which allowed for a comprehensive assessment of their quality and relevance. A significant portion of these images were manually labeled, providing a valuable set of annotations that could be directly incorporated into the training process.

The manual labeling of the generated images played a crucial role in refining the model. By including these hand-labeled images in the training dataset, the model was exposed to a broader range of examples and variations, which contributed to a more robust learning process. This incorporation of synthetic images, along with their accurate labels, led to a notable improvement in the model's accuracy. The model's performance metrics, including correlation and overall accuracy, showed substantial gains, indicating that the additional data had a positive impact.

The results of this process provide compelling evidence that synthetic images can be an effective tool for data augmentation. By generating and incorporating synthetic images into the training dataset, it is possible to enhance the model's ability to generalize and perform more accurately. The improved accuracy demonstrates that these augmented data points add valuable information to the training process, allowing the model to learn from a more diverse and comprehensive set of examples.

This approach highlights the potential of generated images to address data limitations and contribute to the efficiency of model training. In scenarios where obtaining and labeling real data is challenging or resource-intensive, synthetic data serves as a viable alternative that can augment the dataset without the need for extensive manual effort. The successful integration of generated images into the training process underscores the importance of data augmentation in developing robust and accurate models.

Overall, the experience confirms that incorporating synthetic images into model training can lead to significant improvements in performance. By leveraging the capabilities of image generation and data augmentation, it is possible to enhance model accuracy and efficiency, demonstrating the practical value of synthetic data in machine learning and artificial intelligence applications.

## 5 Discussion

Initially, the model's performance in terms of correlation and accuracy was quite low. This indicated that the model was struggling to effectively learn from the existing data, leading to unsatisfactory predictions and a need for enhancement. To address this, a significant effort was made to generate a large volume of synthetic images—specifically, 10,000 images were created.

These generated images were then subjected to predictions, which allowed for an evaluation of their quality and utility. Out of these, many images were manually labeled to provide accurate annotations. The inclusion of these hand-labeled images into the training dataset was a crucial step in improving the model's performance.

By integrating these labeled synthetic images into the training process, the model was exposed to a wider range of examples and variations. This exposure enriched the training data, allowing the model to learn from a more diverse set of inputs. The result of this integration was a noticeable improvement in the model's accuracy. The metrics for correlation and overall performance showed significant gains, indicating that the synthetic images were beneficial.

This outcome demonstrates that synthetic images can be an effective means of data augmentation. Generating and incorporating such images into the training set helped address the limitations of the original dataset, leading to more accurate and robust model performance. The improvement in accuracy underscores the value of using synthetic data to enhance the training process, especially when dealing with challenges related to data availability or labeling.

The experience confirms that synthetic data can play a vital role in model development. By leveraging generated images for data augmentation, it is possible to enhance the efficiency and effectiveness of model training. This approach highlights the practical benefits of synthetic data in machine learning, showing that it can significantly contribute to achieving better model performance and accuracy.

## 6 Conclusion

Initially, the model exhibited low correlation and accuracy, reflecting its struggle to effectively learn from the existing dataset. To tackle this issue, a substantial step was taken to generate a large volume of synthetic images—specifically, 10,000 images were produced.

These synthetic images were subjected to predictions to evaluate their quality and usefulness. Subsequently, many of these images were meticulously hand-labeled, providing valuable ground truth data that could be used for training purposes. The incorporation of these labeled synthetic images into the training process marked a significant advancement.

By adding the hand-labeled synthetic images to the training set, the model benefited from a broader range of data. This additional data enriched the training process by introducing more variability and examples, which in turn enhanced the model's learning capabilities. The model's performance metrics, including accuracy and correlation, improved considerably as a result of this enriched dataset.

The successful improvement in accuracy demonstrates the effectiveness of using synthetic images for data augmentation. The ability to generate and label synthetic images provided a practical solution to the challenge of limited data, contributing valuable information that significantly enhanced the model's performance. This approach underscores the potential of synthetic data to play a crucial role in improving model training and accuracy.

Overall, the results affirm that synthetic images are a valuable tool for augmenting datasets. By generating and integrating these images into the training process, it is possible to achieve more accurate and robust model performance, showcasing the practical benefits of synthetic data in enhancing machine learning outcomes.

## References

[1] "Fine-tuning Kandinsky-2 for Text-to-Image Synthesis," *GitHub*, [Online]. Available: https://github.com/huggingface/diffusers/blob/main/examples/kandinsky2_2/text_to_image/README.md#kandinsky22-text-to-image-fine-tuning.

[2] "Stable Diffusion V1-5 for Diverse Image Synthesis," *Hugging Face Models*, [Online]. Available: https://huggingface.co/runwayml/stable-diffusion-v1-5.

[3] "Exploring Diffusion Models for Fine-grained Image Synthesis," *GitHub*, [Online]. Available: https://github.com/hojonathanho/diffusion.

[4] "SDXL-Lightning: Enhancing Stable Diffusion with Cross-Lingual Learning," *Hugging Face Models*, [Online]. Available: https://huggingface.co/ByteDance/SDXL-Lightning.

[5] Tan, M., and Le, Q. V., "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," *arXiv*, 2019. [Online]. Available: https://arxiv.org/abs/1905.11946.

[6] Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S., "GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium," *arXiv*, 2017. [Online]. Available: https://arxiv.org/abs/1706.08500.

[7] NVIDIA, "NVIDIA A100 Tensor Core GPU Architecture," [Online]. Available: https://www.nvidia.com/en-us/data-center/a100/.

[8] Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R., "SDXL: Enhancing Self-Distillation with Cross-Lingual Learning Capabilities," *arXiv*, 2023. [Online]. Available: https://arxiv.org/abs/2307.01952.

[9] "Image Generation Independent Study," GitHub, [Online]. Available: https://github.com/kbharat7/ImageGen_IndependentStudy

[10] Lakhani, Paras and Sundaram, Baskaran, "Deep learning at chest radiography: Automated classification of pulmonary tuberculosis by using convolutional neural networks," *Radiology*, vol. 284, no. 2, pp. 574-582, 2017.

[11] Rajpurkar, Pranav and others, "CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning," *Stanford Machine Learning Group*, 2017. [Online]. Available: https://arxiv.org/abs/1711.05225.