# A Project Report on Image Generation Using Stable Diffusion and Control Net

## Thought Process

The primary goal was to develop an efficient image generation system that uses both text prompts and depth maps for added control over the output images. We used **ControlNet** with Stable Diffusion to allow better conditioning of the generated images. Created a pipeline for successful execution

Main Considerations:

- **Efficiency**: Using a pre-trained model saved locally to reduce the time and cost of downloading from external sources.
- **Scalability**: Writing the code in such a way that allows easy scaling for generating multiple images based on various prompts and conditions.
- **Latency Measurement**: Added a functionality to calculate and display the time taken for image generation, which is crucial for performance analysis.

## Visual Results

The generated images were saved in the generated_images/ folder(present in the github repo).
Each image corresponds to a prompt and depth map, allowing for both textual and visualization of the output.

Sample images results:

- **Prompt**: "A luxury bedroom,"
  - o Generated Image: A well generated image of a bedroom, showing details like a bed and lights as seen when we google luxury bedroom on google .
  - o Edge Map: Canny edge detection applied to the depth map provided an accurate boundary for the image.
- **Aspect Ratio Analysis**: Images generated with different aspect ratios showed that 1:1 ratio did not provide with the best quality as the image got short and elements were overlapped, while wider ratios like 16:9 led to some loss of image on the corner/edges of images but helped in identifying the features correctly.

# Analysis of Performance

## a. Areas where it will work well:

1. <u>Image Conditioning</u>: Using depth maps and edge detection with ControlNet provides highly detailed images. The additional conditioning ensures that the generated images closely follow the depth information, adding realism to the images.
2. <u>Pre-trained Model Use</u>: Loading the pre-trained model from a local directory significantly reduced the time required to initialize the pipeline.
3. <u>Latency Monitoring</u>: The addition of latency measurement helped assess how long each image took to generate, allowing us to monitor performance and look for optimizations. <u>On an average each image took around 35 minutes to generate.</u>
4. <u>Edge Detection</u>: Using Canny edge detection with depth maps adds an extra layer of control to the generated images, ensuring that key features are retained.

## b. Areas where it will fail:

1. <u>Handling Large Depth Maps</u>: For large depth maps, the process becomes slow, and in some cases, the images fail to generate. The current approach does not handle memory-intensive operations efficiently.
2. <u>Aspect Ratio Distortions</u>: Non-standard aspect ratios lead to distortion in some image areas. While we attempted to preserve image quality, there are visible artifacts in wide or narrow aspect ratio.
3. <u>Reduced Latency:</u> When the latency of an image is less , it will lead to less accurate response and sometimes no image generation will take place.

# Ideas for Improvement

1. <u>Memory Optimization</u>: One idea is to resize depth maps dynamically before processing, reducing the load on the GPU.
2. <u>Better Aspect Ratio Handling</u>: Implement padding techniques to maintain image quality across all aspect ratios without distortion.
3. <u>Edge Detection Enhancement</u>: Explore using more advanced edge detection techniques (e.g., Sobel, Laplacian) that may work better with different types of depth maps, especially when the depth data is noisy or incomplete.
4. <u>Batch Processing</u>: Implement a batch processing pipeline to generate images more efficiently for multiple prompts, reducing the time taken by individual function calls
5. <u>Latency Improvement:</u> We can significantly improve the latency of the images upto 20-50% by using limiting the output length, caching frequently and saving the responses in the memory so that if the image is generated again it will be fater.
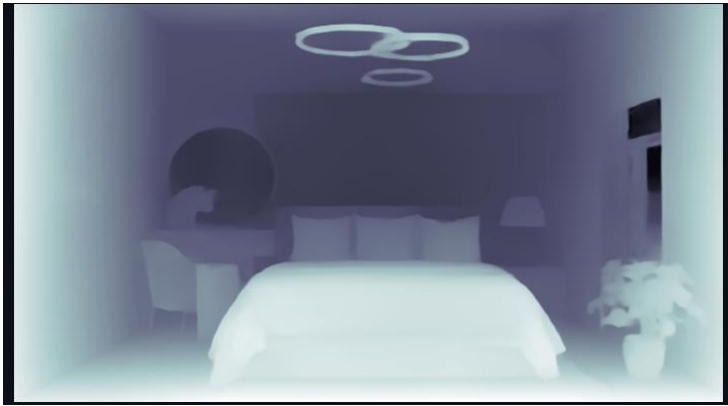
## Some Output Images



Downloading the model from hugging face and extracting pipeline components.

```
(venv) PS D:\New folder\avataar_dummy_ass_!> pip install opencv
k\pytorch\pytorch\builder\windows\pytorch\aten\src\ATen\native\transformers\cuda\sdp_utils.cpp:555.)
    attn_output = torch.nn.functional.scaled_dot_product_attention(
100%|                                                              | 50/50 [43:03<00:00, 51.68s/it]
12%|                                    14%|                        | 50/50 [1:04:52<00:00, 77.85s/it]
100%|                                                              | 50/50 [41:09<00:00, 49.38s/it]
100%|                                                              | 50/50 [40:49<00:00, 48.99s/it]
100%|                                                              | 50/50 [41:35<00:00, 49.90s/it]
100%|                                                              | 50/50 [41:09<00:00, 49.40s/it]
58%|                                    | 29/50 [47:16<3 60%|
100%|                                                              | 50/50 [1:30:34<00:00, 108.69s/it]
```

Model Training and images being generated



Aspect ratio (1:1)



Aspect Ratio (16:9)