

From Pixels to Pictures: Understanding the Internal Representation of Latent Diffusion Models

Karina Chen
kac009@ucsd.edu

Atharva Kulkarni
apkulkarni@ucsd.edu

Ester Tsai
etsai@ucsd.edu

Zelong Wang
zew013@ucsd.edu

Mentor: Alex Cloninger
acloninger@ucsd.edu

Mentor: Rayan Saab
rsaab@ucsd.edu



Project Background

What is Stable Diffusion?

- Stable Diffusion is an open-source diffusion model that generates images from text prompts.
- Stable Diffusion is a two-stage framework that consists of:
 - A latent diffusion model (LDM)
 - The LDM learns to predict and remove noise in the latent space by reversing a forward diffusion process.
 - A variational autoencoder (VAE)
 - The VAE converts data between latent and image space.
 - After the LDM synthesizes a denoised latent z , the decoder of VAE converts the denoised latent z to the image space.

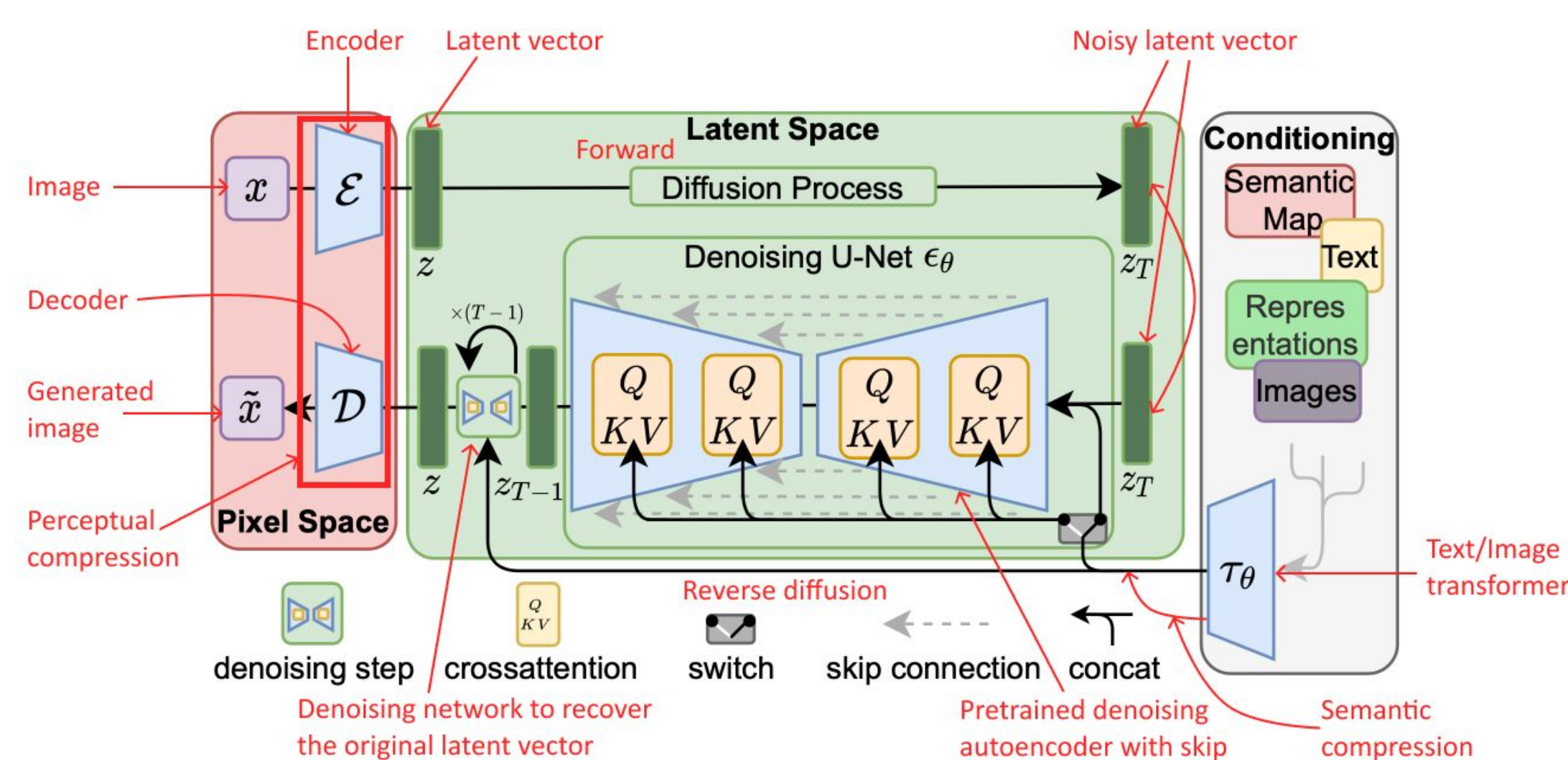


Figure 1. Architecture of an LDM (Rombach, Blattmann 2022)

Problem Statement

- Does an LDM create an internal 3D representation of the object it portrays?
- How early in the denoising process do depth, saliency, and shading information develop in the internal representation?
- At what time step does an image classifier correctly detect the object?

Data

617 images (512 pixels x 512 pixels) generated using Stable Diffusion v1.4



Image generated by Stable Diffusion v1.4 using the text prompt "ZIGGY - EASY ARMCHAIR" and seed 64140790.



Salient object detection mask generated by TRACER.



Shading and illumination map generated by Intrinsic.



Depth map generated by MiDaS.

Figure 2. Ground truth images

Internal Representation

Methods

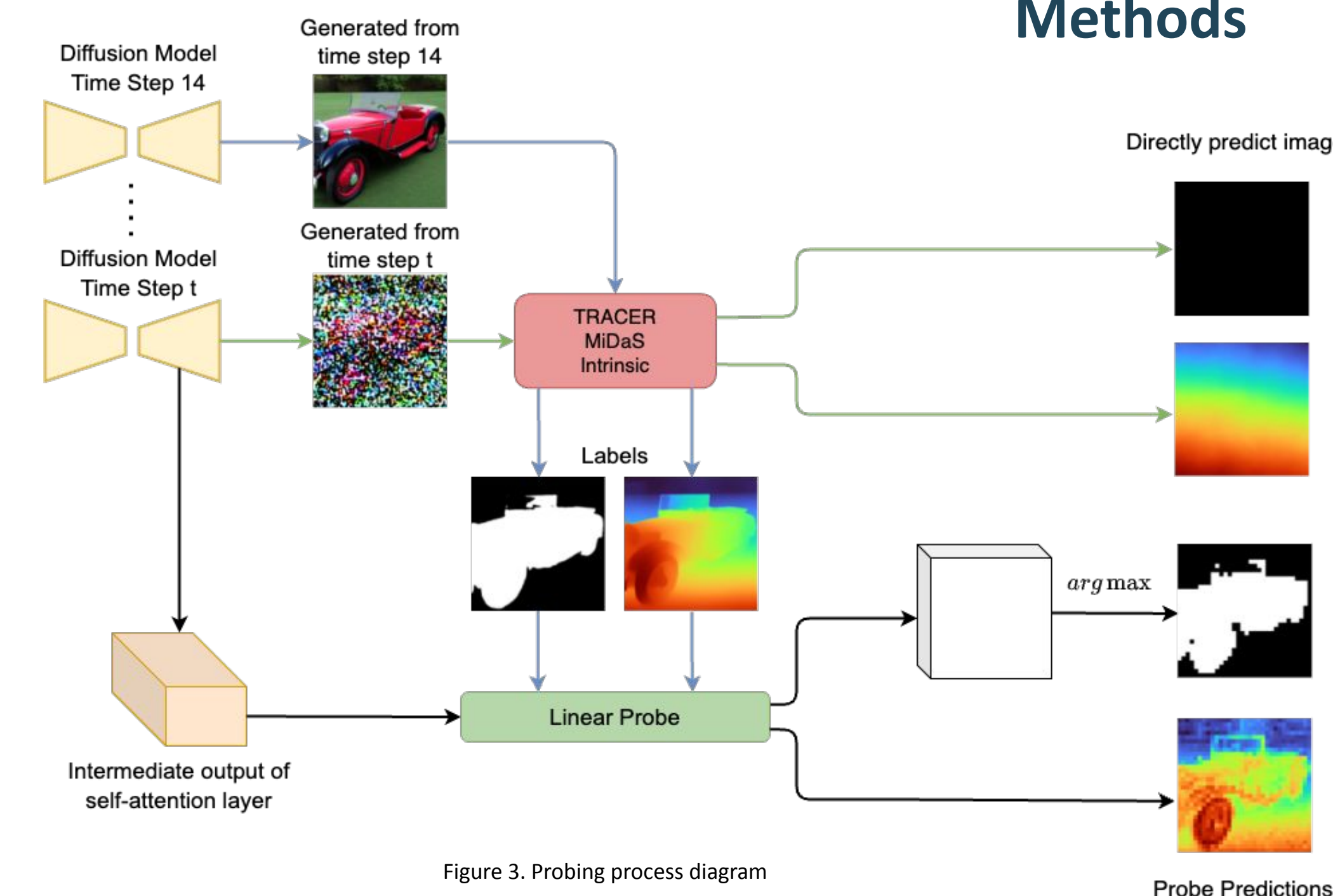


Figure 3. Probing process diagram

Results: Probing the LDM

Probe performance at the last step	Score between -1 and 1
Foreground Segmentation Dice Coefficient	0.85
Depth Estimation Rank Correlation	0.71
Shading Estimation Rank Correlation	0.62

- Using intermediate activations of noisy input images, linear probes can accurately predict the foreground, depth, and shading.
 - Shown by high Dice Coefficient and Rank Correlation in the table.
- All three properties emerge early in the denoising process (around step 3 out of 15), suggesting that the spatial layout of the generated image is determined at the very beginning of the generative process.

3D properties in LDM emerge at step 3

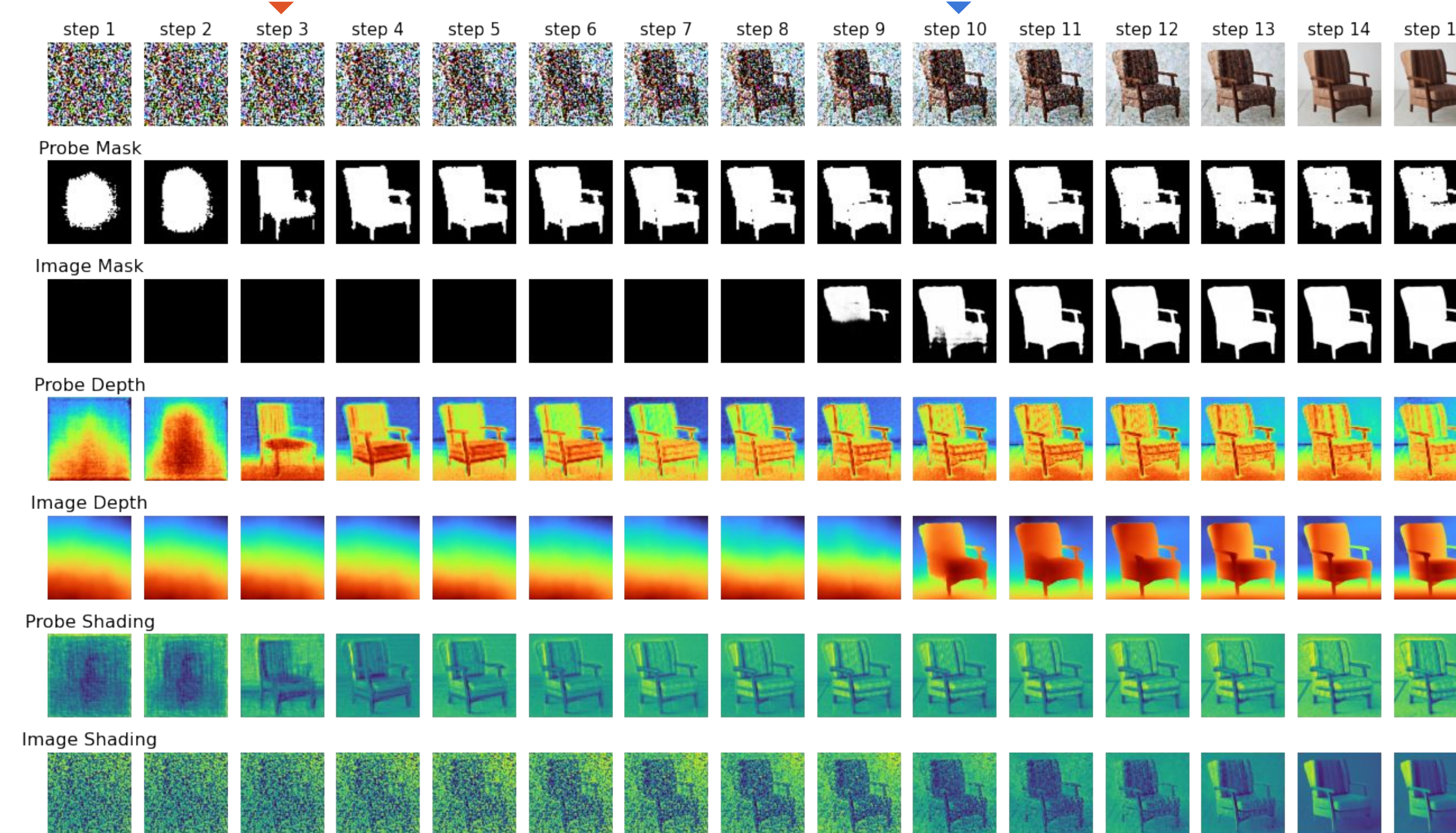


Figure 4. Intermediate steps for the generated image, probe, and model results

Image Classification

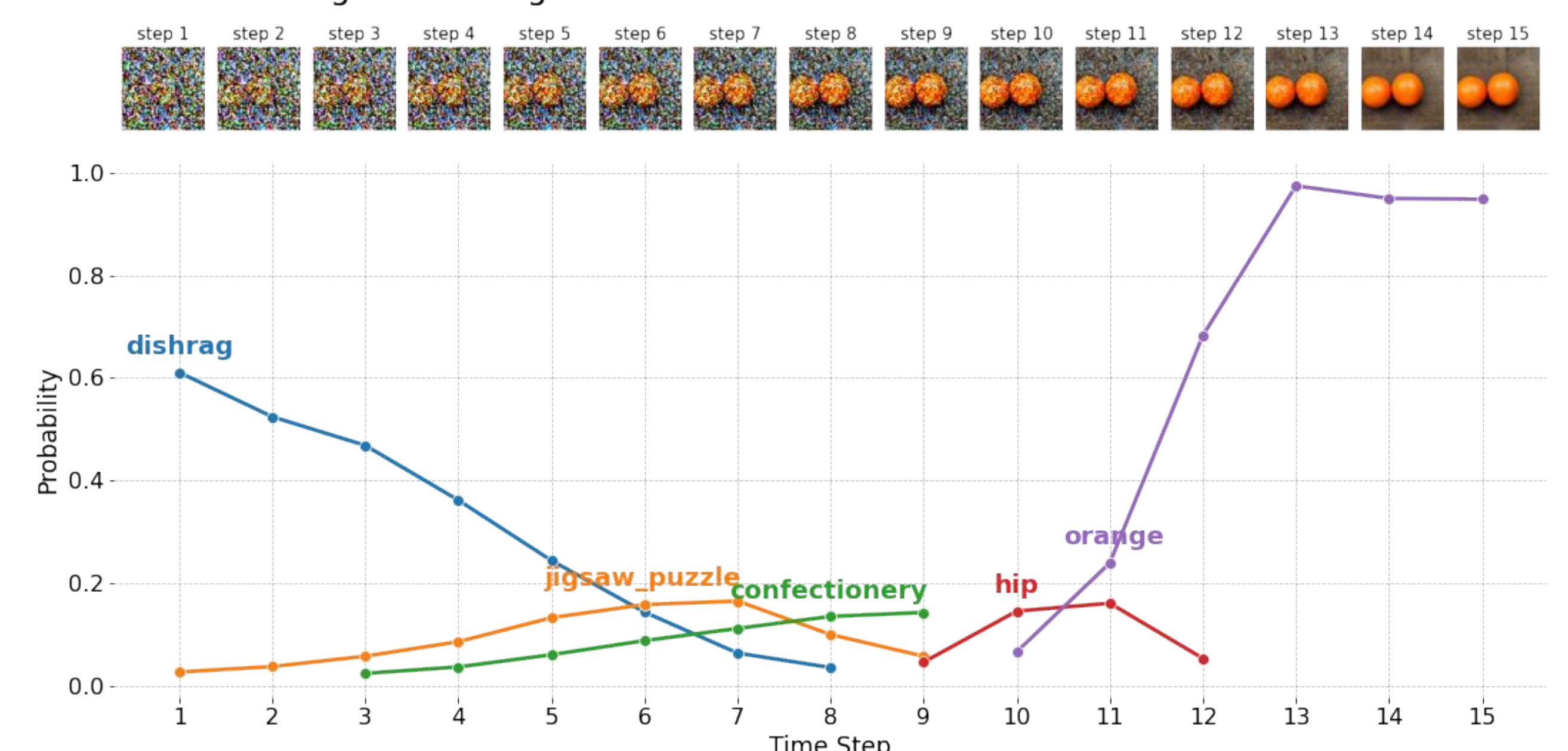
Methods

- Generate images using Stable Diffusion with prompts that match ImageNet categories.
 - For example, prompt = "lemon".
- Run each intermediate image through VGG-16, an image classification model trained on ImageNet.
- Visualize predictions results.

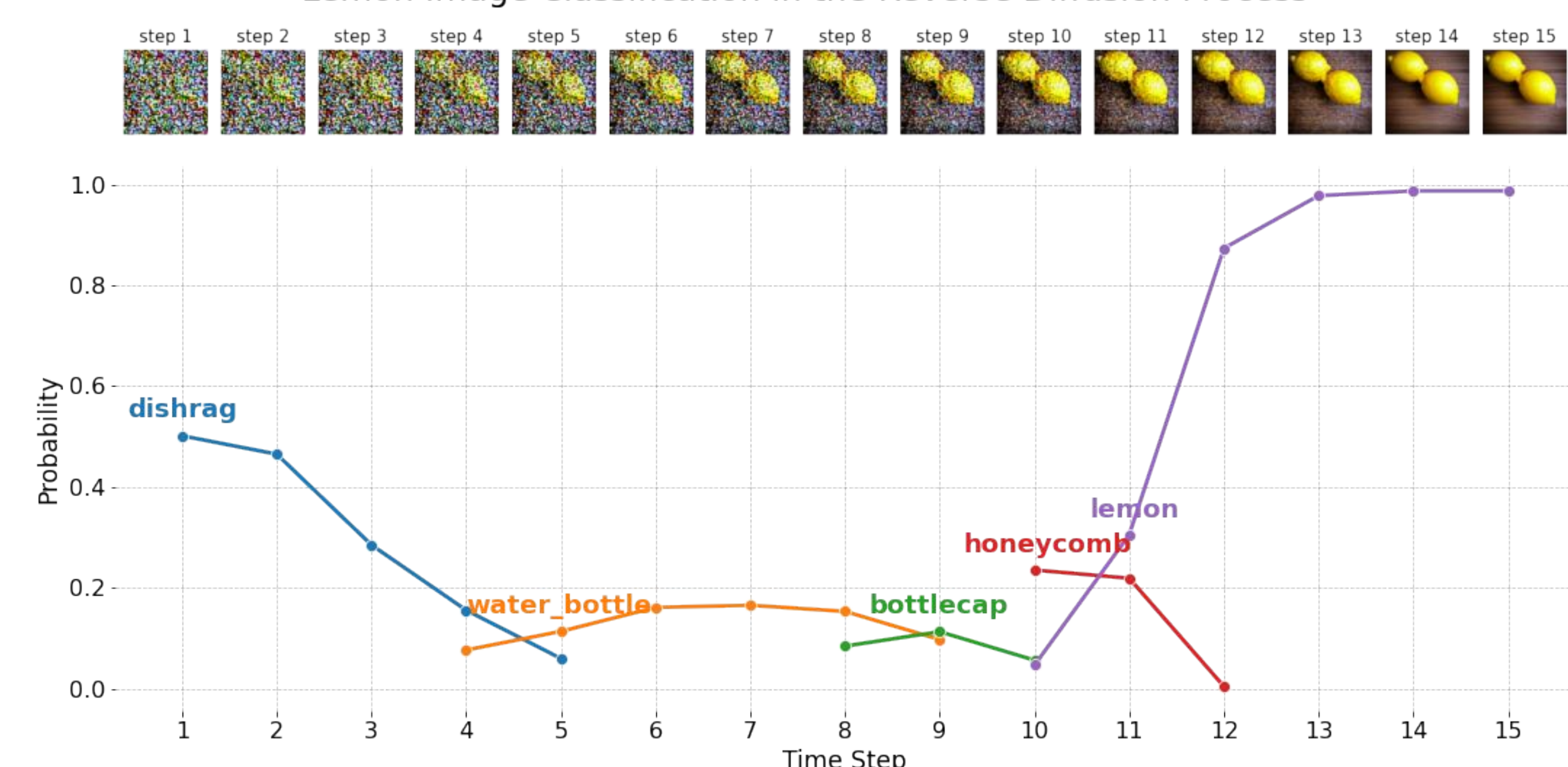
Results

- Comparing classification confidence for generated vs. real images.
 - Generated images: two lemons (98.75%), two oranges (94.8%).
 - Real images: two lemons (99.4%), singular lemon (87.7%), singular orange (87.0%).
- The correct classification has high confidence (> 90%) towards the end of the diffusion process for the majority of generated images.
 - This means that the generated images are fairly good representations of the object prompted.

Orange Fruit Image Classification in the Reverse Diffusion Process



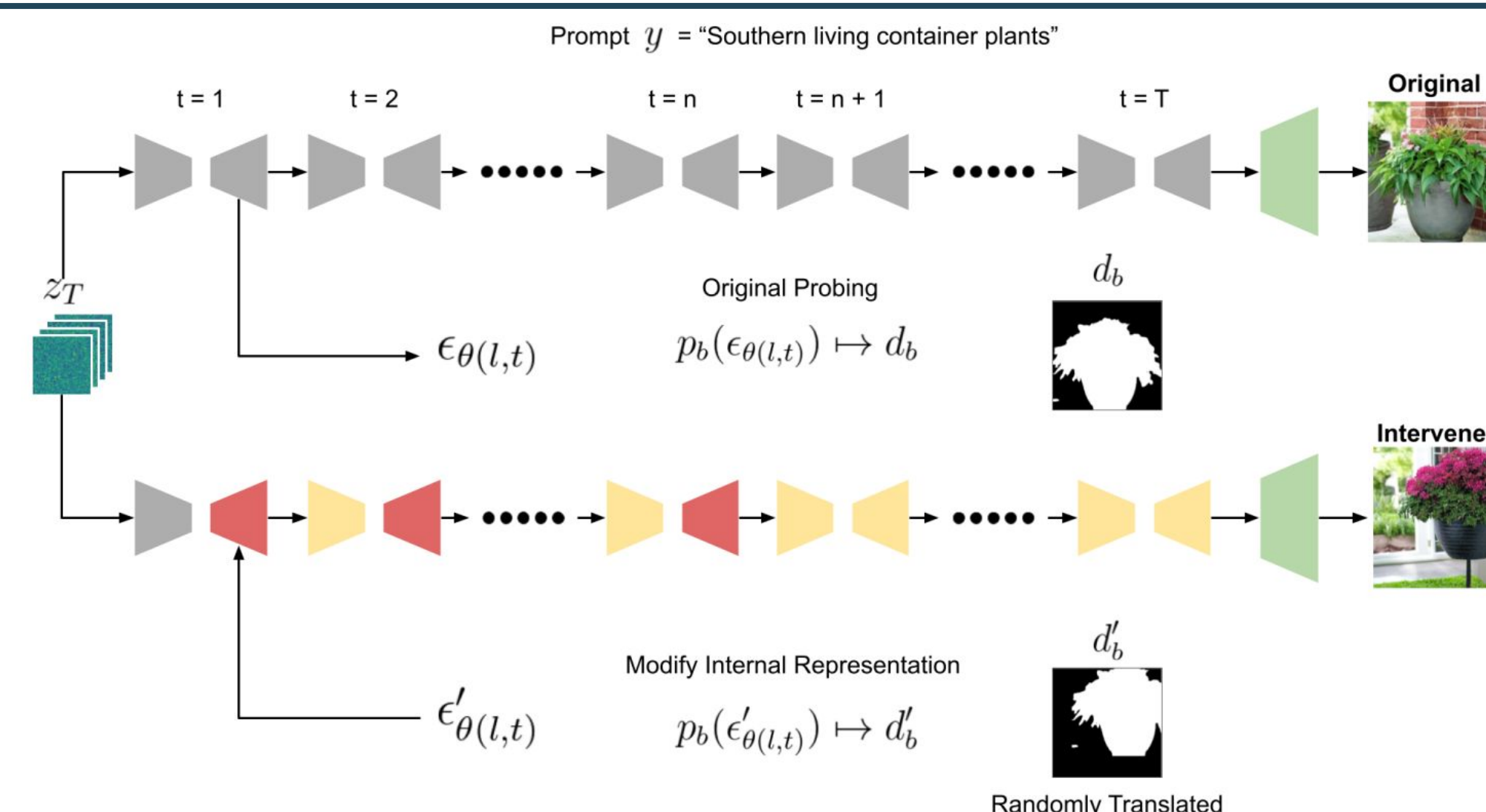
Lemon Image Classification in the Reverse Diffusion Process



Future Works: Intervening the LDM

Figure 5. The Intervention workflow (Chen, 2023).

- The foreground object can be repositioned by modifying the activations of the U-Net decoders.
- First, obtain a target mask by translating the original mask.
 - Goal: to find the activations (i.e. probe inputs) that cause the probe to output a mask highly similar to the target mask.
- Perform gradient descent on the activations until the probe can output the desired target mask.
- Replace the original activations with the modified activations, then resume the denoising process.



- Foreground mask has a causal role in image generation.
- Intervention: Without changing the prompt, input latent vector, and model weights, we can modify the scene layout of generated image by editing the foreground mask (Y. Chen et al.).

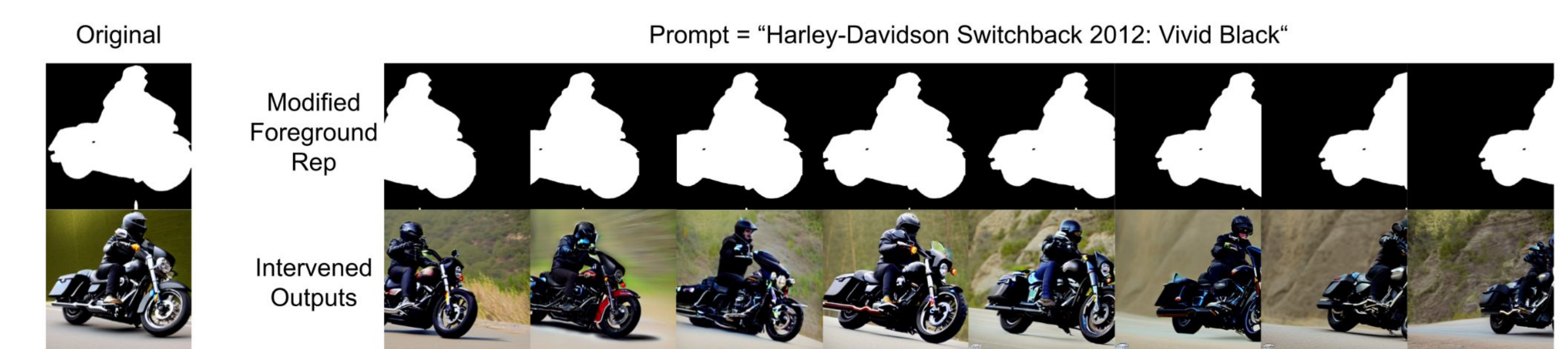


Figure 6. Intervening the LDM to produce different outputs (Chen, 2023)

**HELLO,
Ester**



From Pixels to Pictures: Understanding the Internal Representation of Latent Diffusion Models

Karina Chen
kac009@ucsd.edu

Atharva Kulkarni
apkulkarni@ucsd.edu

Ester Tsai
etsai@ucsd.edu

Zelong Wang
zew013@ucsd.edu

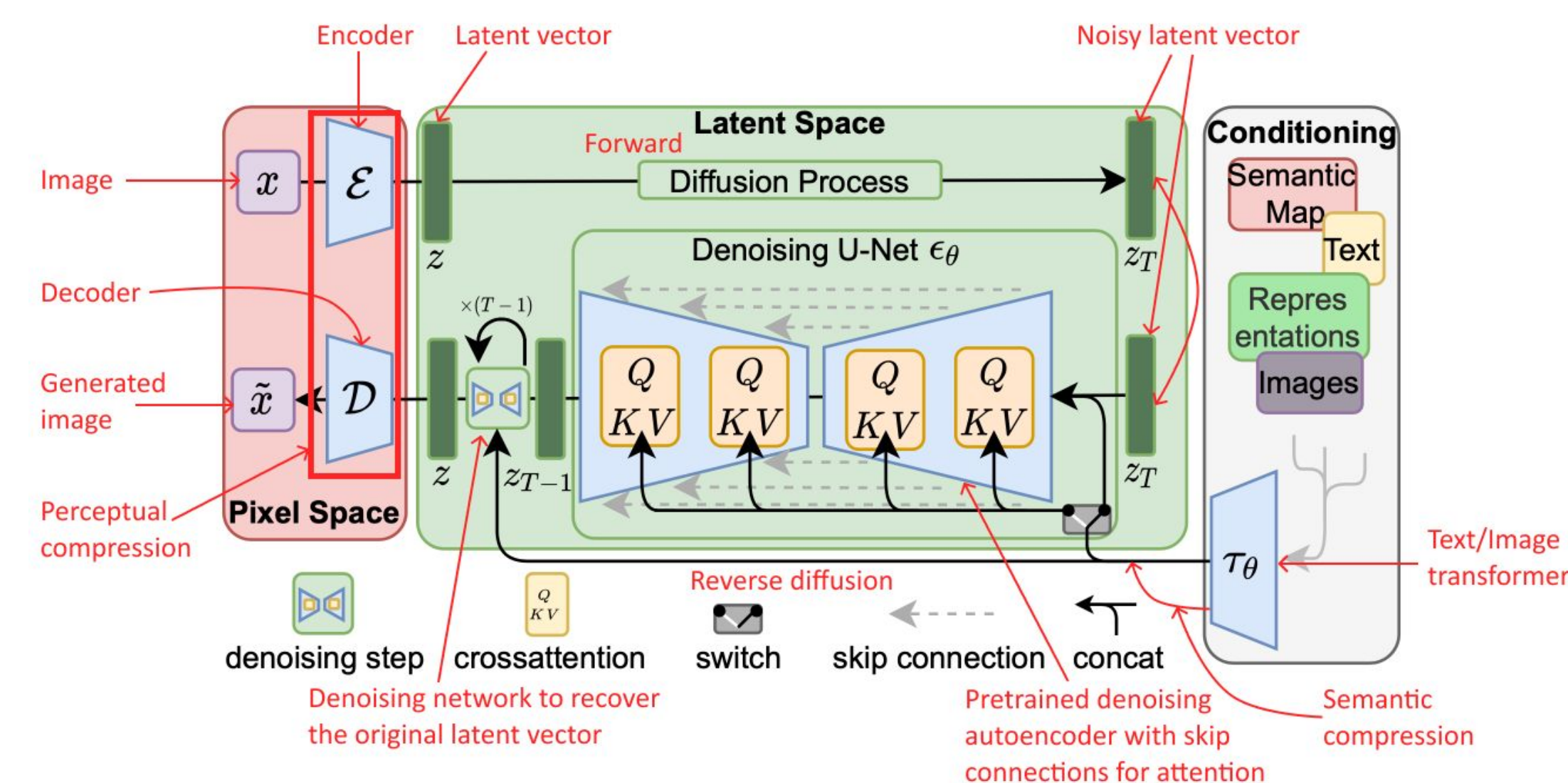
Alex Cloninger
acloninger@ucsd.edu

Rayan Saab
rsaab@ucsd.edu

Project Background

What is Stable Diffusion?

- Stable Diffusion is an open-source diffusion model that generates images from text prompts.
- Stable Diffusion is a two-stage framework that consists of:
 - A latent diffusion model (LDM)
 - The LDM learns to predict and remove noise in the latent space by reversing a forward diffusion process.
 - A variational autoencoder (VAE)
 - The VAE converts data between latent and image space.
 - After the LDM synthesizes a denoised latent z , the decoder of VAE converts the denoised latent z to the image space.



Problem Statement

Does an LDM create an internal representation of 3D information?
At what time step does an image classifier correctly detect the object?

Data

617 images (512 pixels x 512 pixels) generated using Stable Diffusion v1.4



Image generated by Stable Diffusion v1.4 using the text prompt "ZIGGY - EASY ARMCHAIR" and seed 64140790.



Salient object detection mask generated by TRACER.



Shading and illumination map generated by Intrinsics.



Depth map generated by MiDaS.

Internal Representation

Methods

- Internal representation** = the neural network's self-attention layer's intermediate activation output
- The linear probe model**
 - Input:** internal representation of a LDM
 - At a specific time step, for a specific block and layer of the U-Net
 - Output:** predicted image showing a certain property (e.g. depth, salient-object detection, shading)

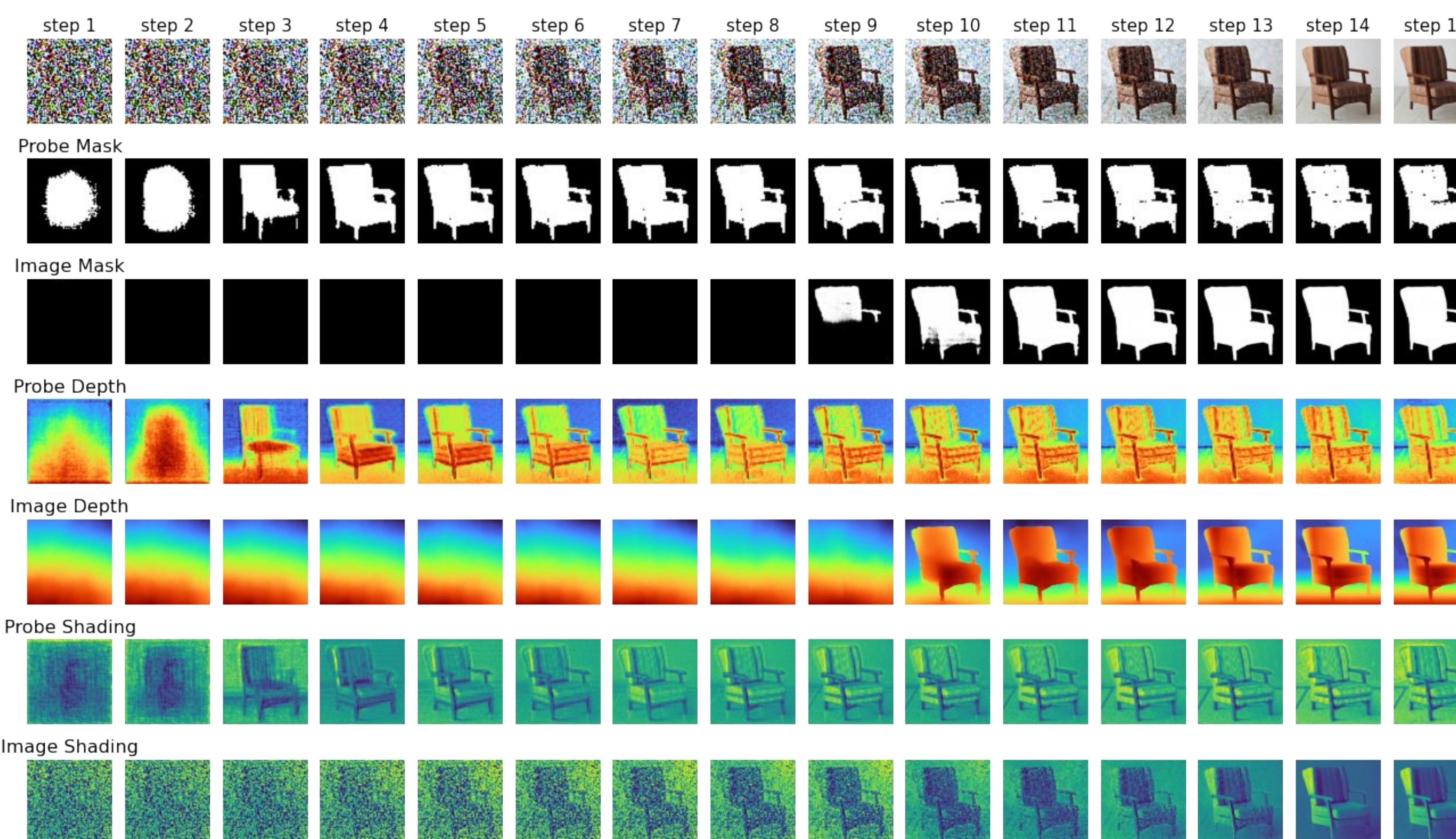
Spatial and Feature Dimensions of Self-Attention Layers in the LDM

Blocks	Number of Self-Attn Layers	Spatial h x w	Feature c
Encoder 1	2	64 x 64	320
Encoder 2	2	32 x 32	640
Encoder 3	2	16 x 16	1280
Encoder 4	0	-	-
Bottleneck	1	8 x 8	1280
Decoder 1	0	-	-
Decoder 2	3	16 x 16	1280
Decoder 3	3	32 x 32	640
Decoder 4	3	64 x 64	320

Results: Probing the LDM

Foreground segmentation Dice coefficient	0.85
Depth Rank Correlation	0.47
Shading Rank Correlation	

- Using intermediate activations of noisy input images, linear probes can accurately predict the foreground, depth, and shading.
- All three properties emerge early in the denoising process (around step 3 out of 15), suggesting that the spatial layout of the generated image is determined at the very beginning of the generative process.



Results: Intervening the LDM

- Foreground mask has a causal role in image generation.
- Intervention: Without changing the prompt, input latent vector, and model weights, we can modify the scene layout of generated image by editing the foreground mask (Y. Chen et al.)

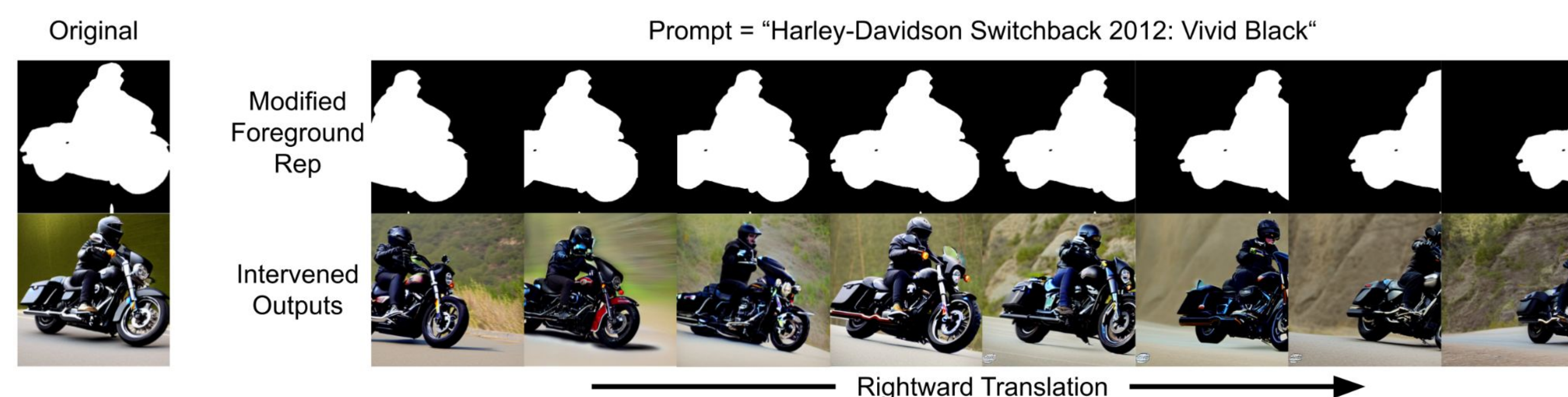


Image Classification

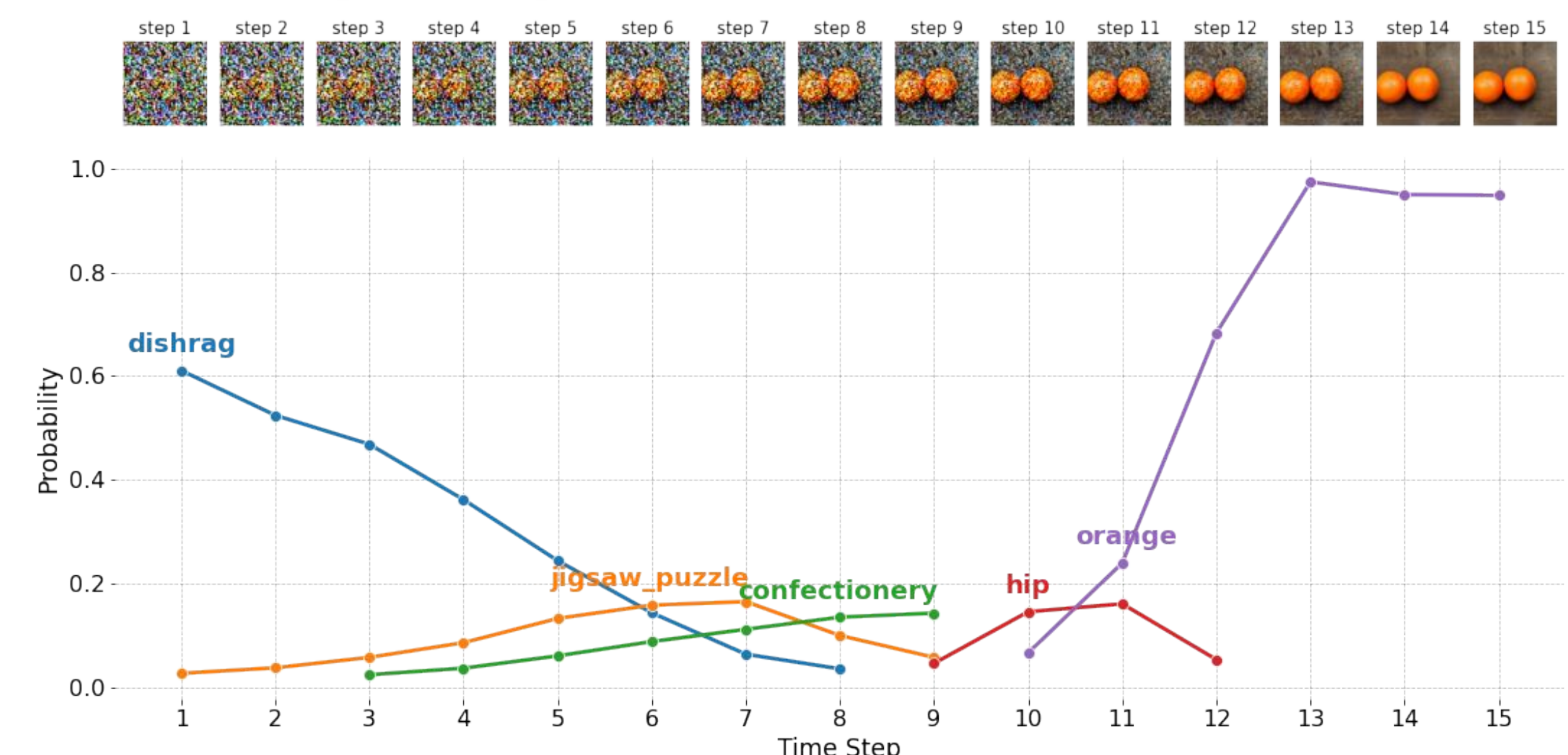
Methods

- Generate images using Stable diffusion with prompts that match ImageNet categories
- Run each intermediate image through VGG-16 (image classification model trained on ImageNet)
- Explore predictions and plot results

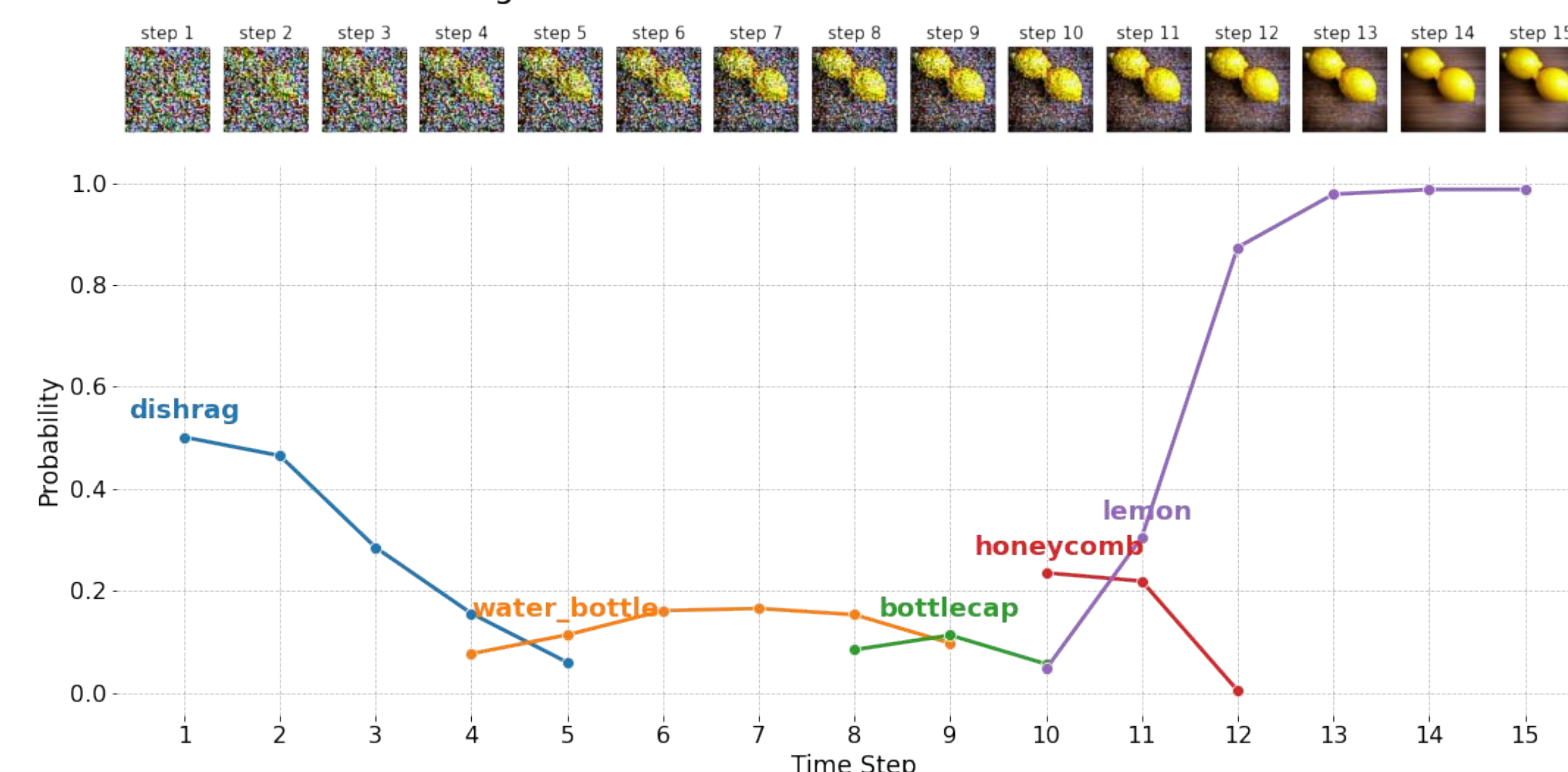
Results

- Comparing classifications for generated versus real images
 - Generated images: two lemons (98.75%), two oranges (94.8%)
 - Real images: singular lemon (87.7%), two lemons (99.4%), singular orange (87.0%)
- The correct classification has high probability (> 90%) towards the end of the diffusion process for the majority of generated images, starting from around step 12.
 - This means that the diffused images are fairly good representations of the object prompted.

Orange Fruit Image Classification in the Reverse Diffusion Process



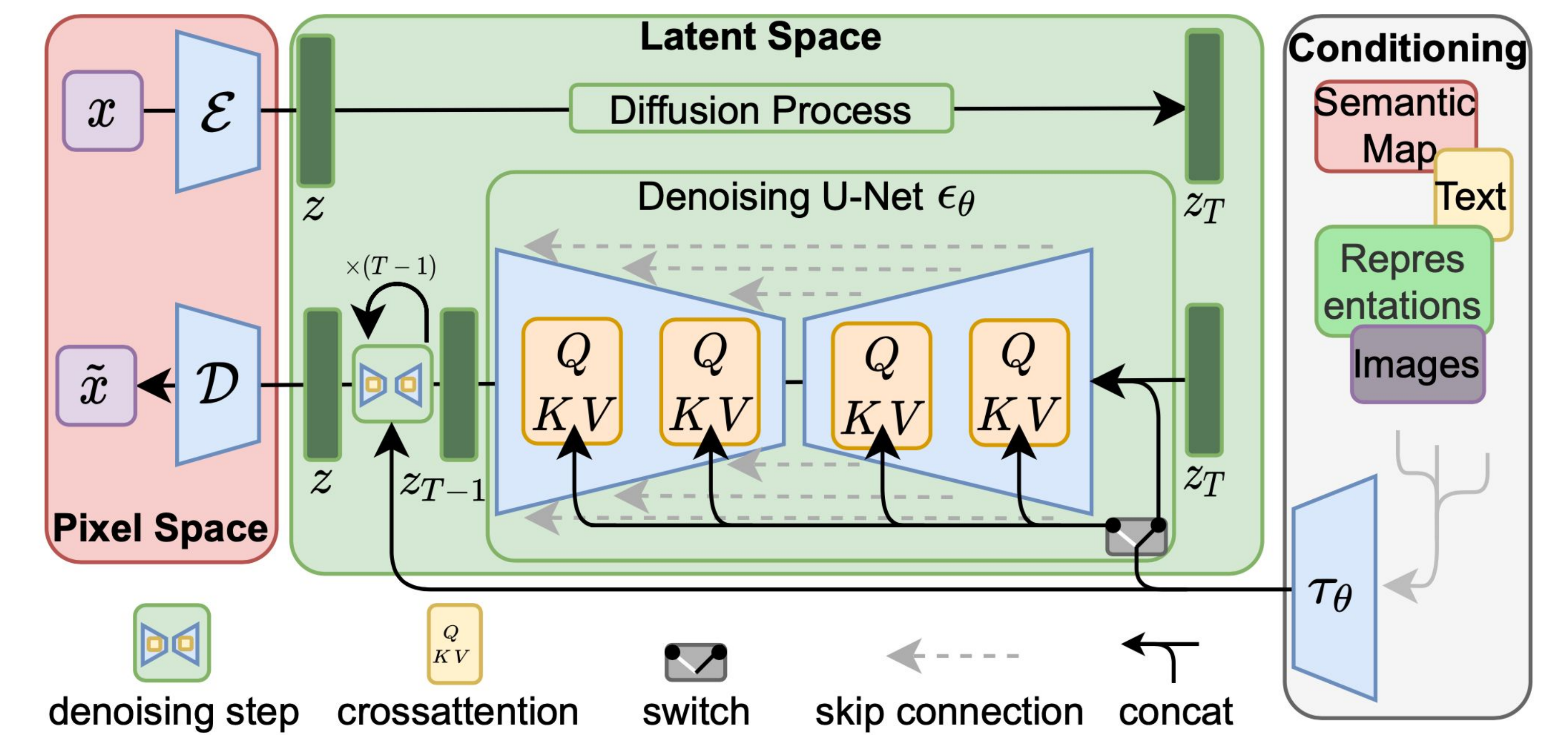
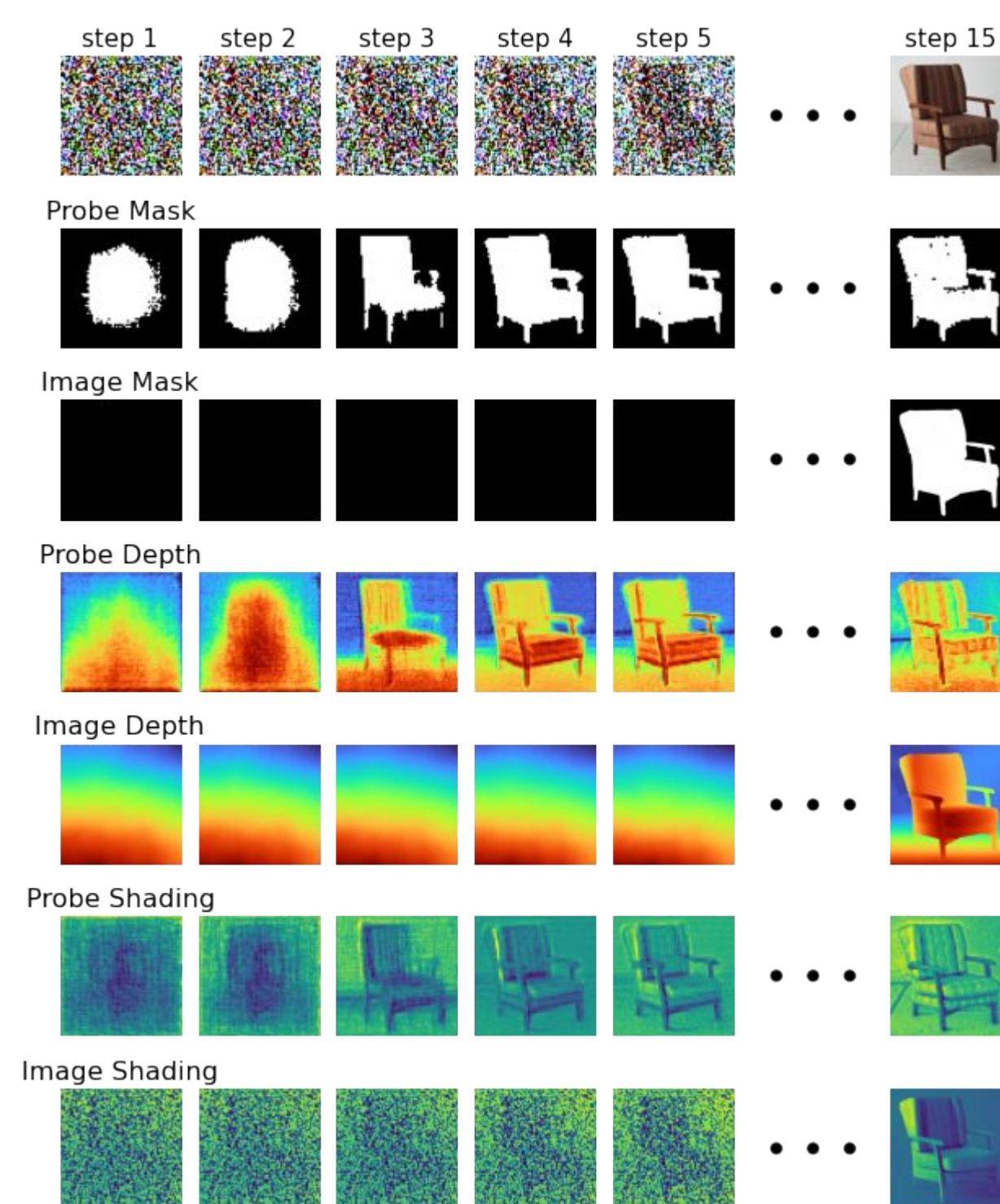
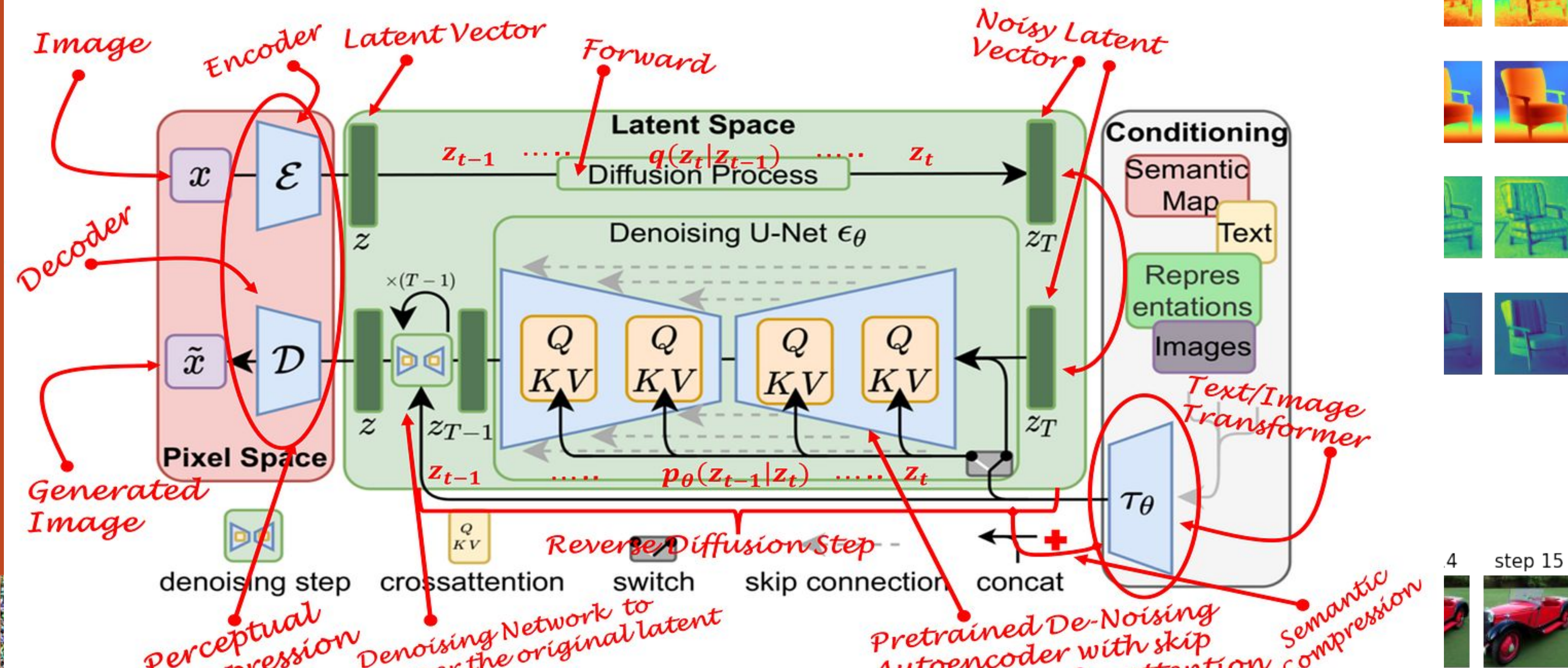
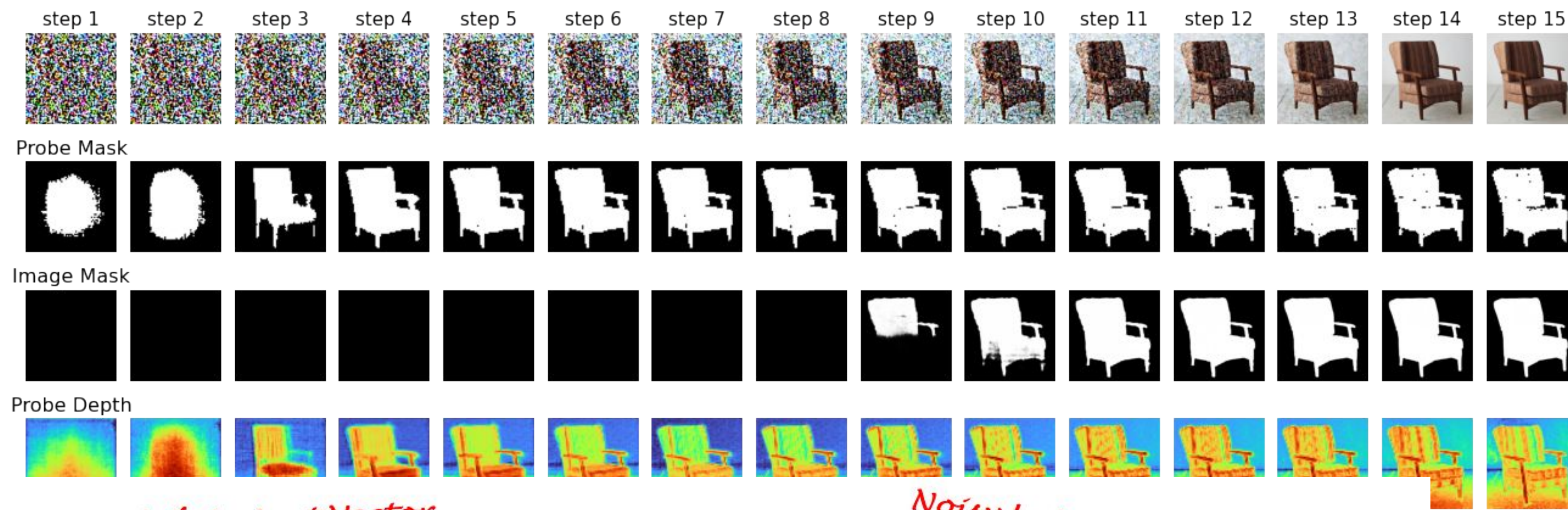
Lemon Image Classification in the Reverse Diffusion Process



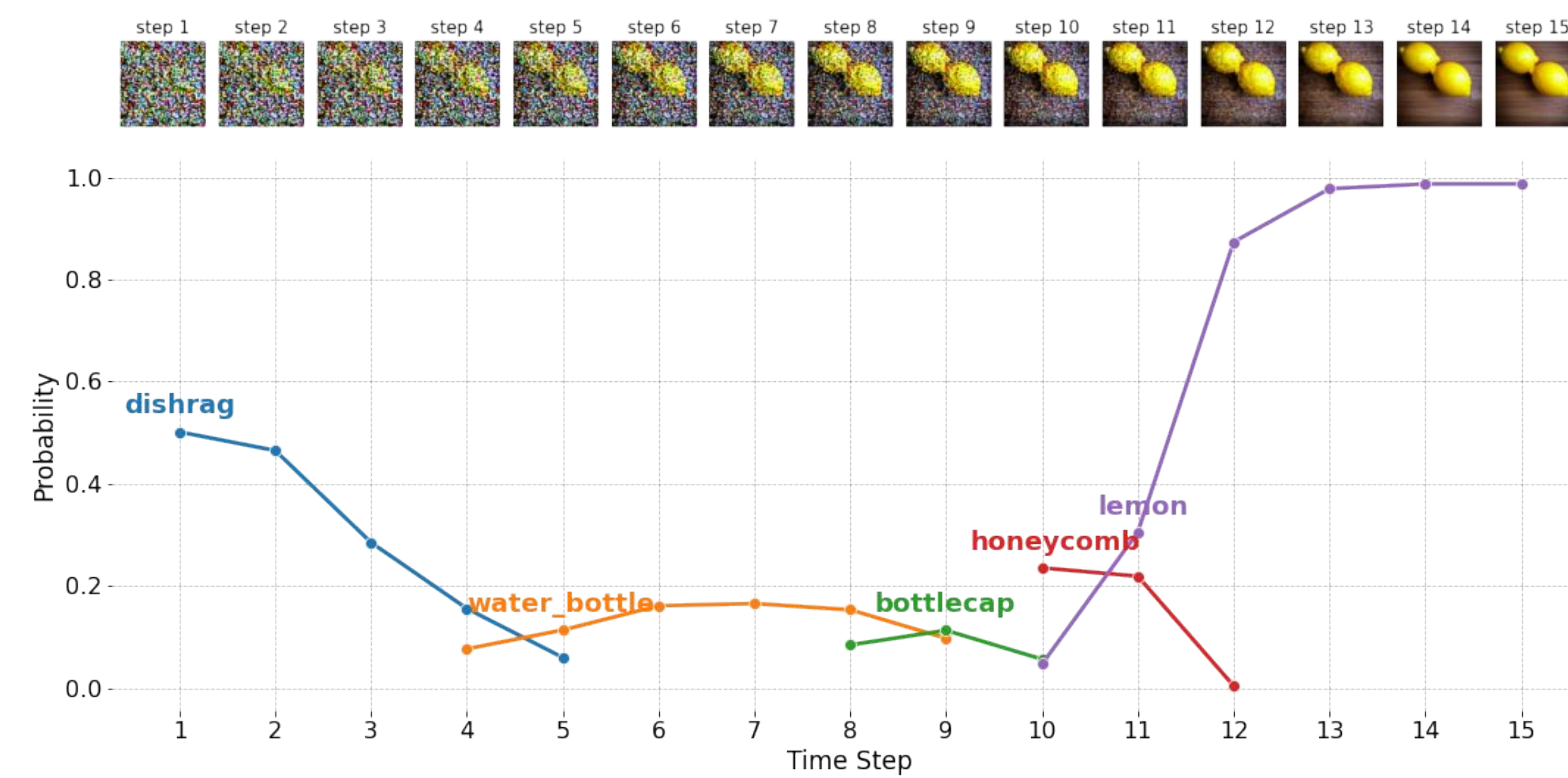
Future Work & Acknowledgements

- Adversarial attacks
 - Explore where diffusion models are deployed, and examine impact of adversarial attack
 - Adversarial training: train diffusion model on clean and adversarially perturbed images to find out robustness
 - Using diffusion models to prevent adversarial attacks
 - Adversarial attacks on diffusion models
- Fairness and bias mitigation
 - Investigate methods for detecting and mitigating biases in generated images (e.g., sensitive attributes such as race, gender, age)
 - Explore techniques for promoting fairness and equity in the outputs of diffusion models, such as adversarial debiasing
- Augmenting datasets
 - Using the encoded depth information to explore augmenting, for example, autonomous vehicle datasets

Our research method draws inspiration from “Beyond Surface Statistics: Scene Representations in a Latent Diffusion Model” (Y. Chen et al.) and moves beyond their focus on depth to explore other image information such as shading and illumination.



Lemon Image Classification in the Reverse Diffusion Process



Future Work

- Adversarial attacks
 - Adversarial training: train diffusion model on clean and adversarially perturbed images to measure robustness against adversarial attacks
- Improve diffusion model training speed
 - If the bulk of the information is already encoded by very early steps in the denoising process, we can potentially speed up the rest of the steps without sacrificing quality
- Fairness and bias mitigation
 - Investigate methods for detecting and mitigating biases in generated images (e.g., sensitive attributes such as race, gender, age)
 - Explore techniques for promoting fairness and equity in the outputs of diffusion models, such as adversarial debiasing
- Augmenting datasets
 - Using the encoded depth information to explore augmenting, for example, autonomous vehicle datasets

Karina Chen
kac009@ucsd.edu

Atharva Kulkarni
apkulkarni@ucsd.edu

Ester Tsai
etsai@ucsd.edu

Zelong Wang
zew013@ucsd.edu

Alex Cloninger
acloninger@ucsd.edu

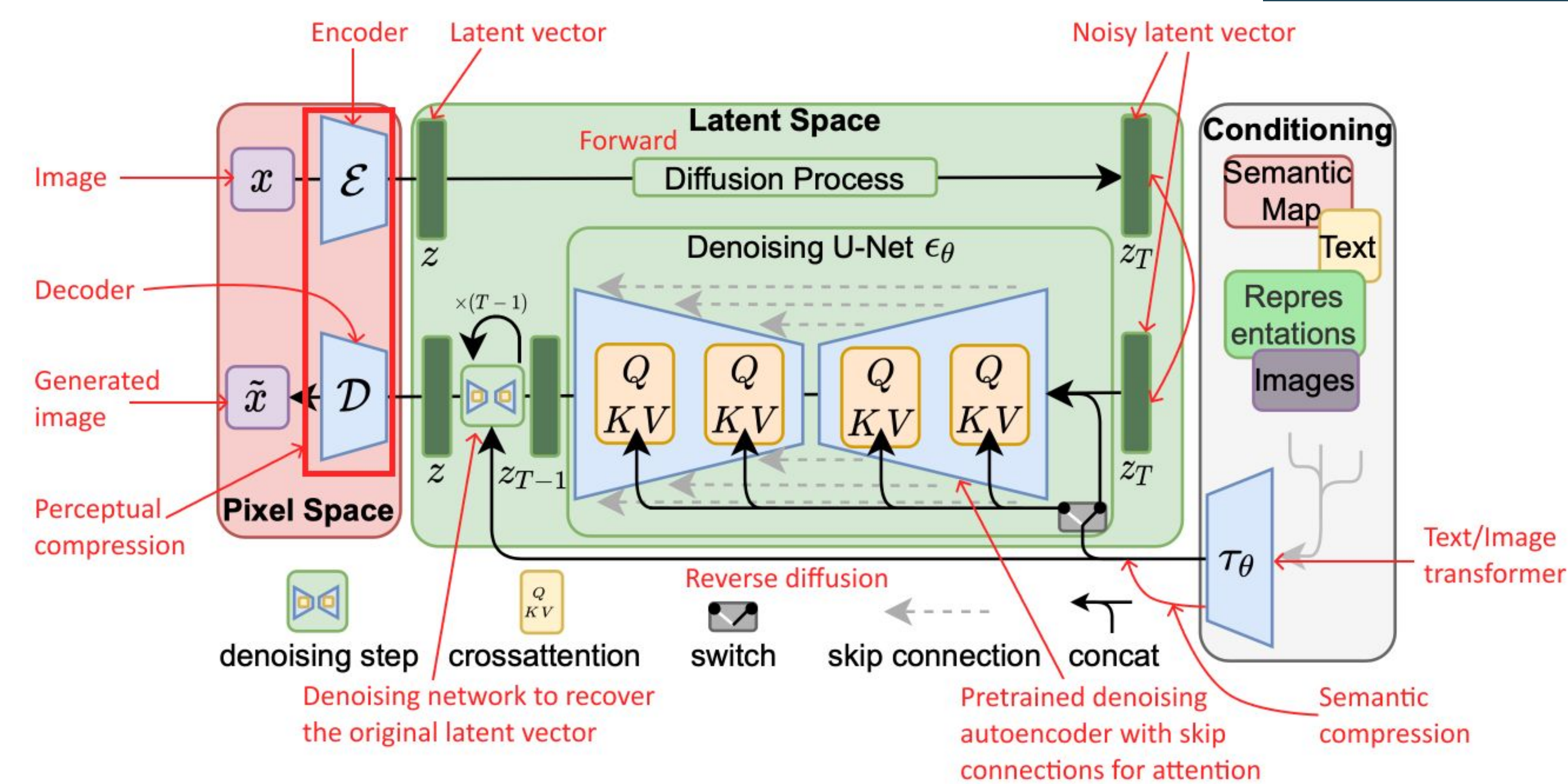
Rayan Saab
rsaab@ucsd.edu

Project Background

Probing the LDM for 3D Information

Image Classification

Results and Discussion



Spatial and Feature Dimensions of Self-Attention Layers in the LDM

Blocks	Number of Self-Attn Layers	Spatial $h \times w$	Feature c
Encoder 1	2	64×64	320
Encoder 2	2	32×32	640
Encoder 3	2	16×16	1280
Encoder 4	0	-	-
Bottleneck	1	8×8	1280
Decoder 1	0	-	-
Decoder 2	3	16×16	1280
Decoder 3	3	32×32	640
Decoder 4	3	64×64	320

Conclusion

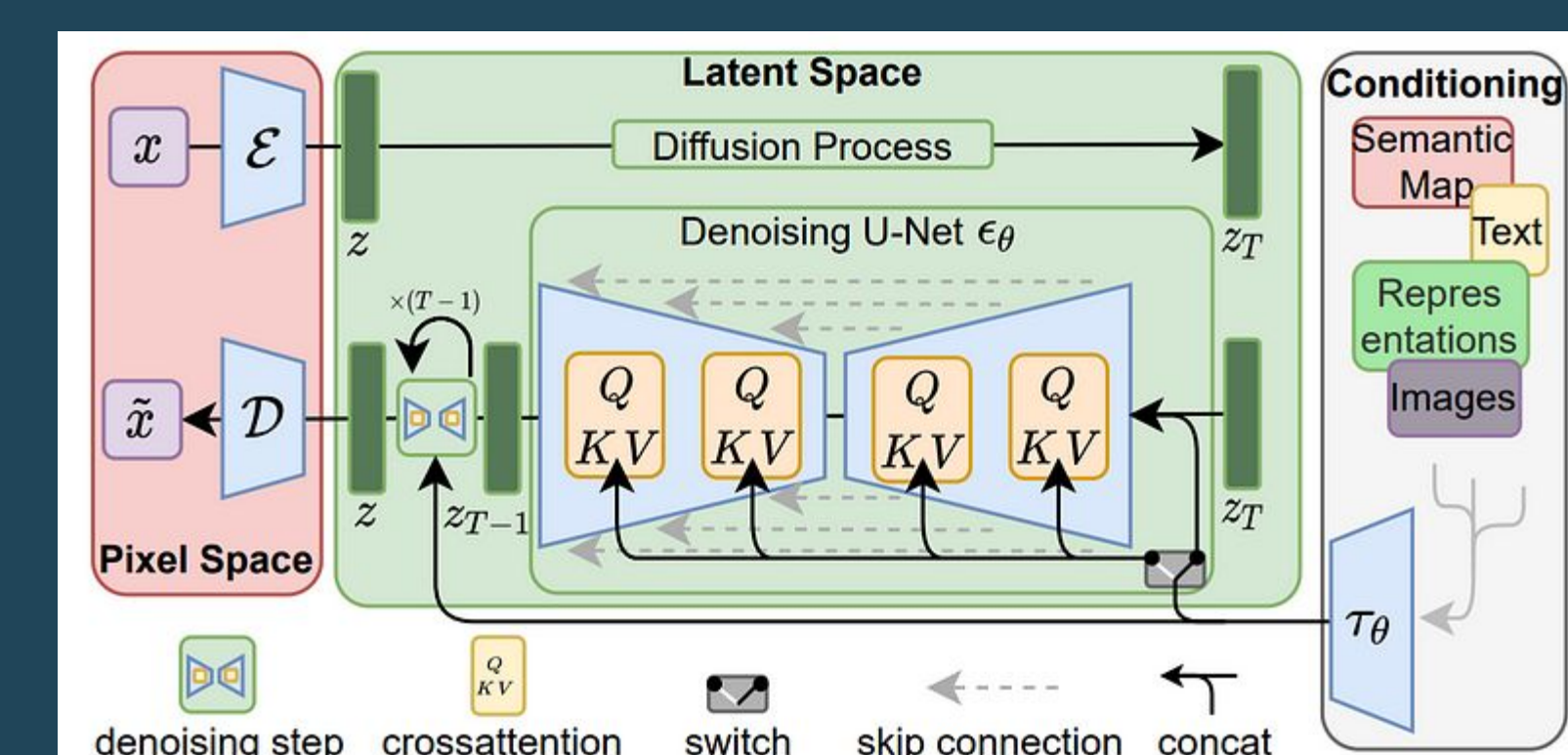
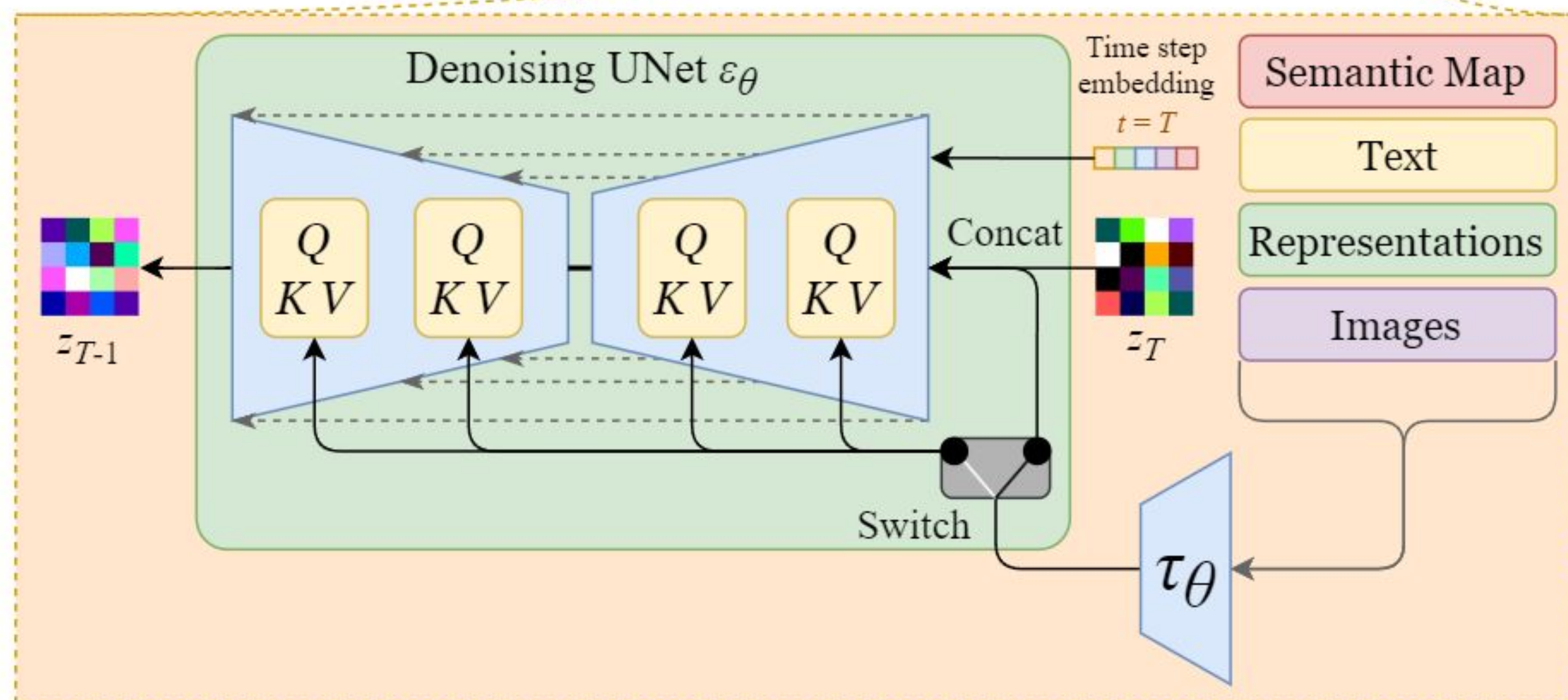
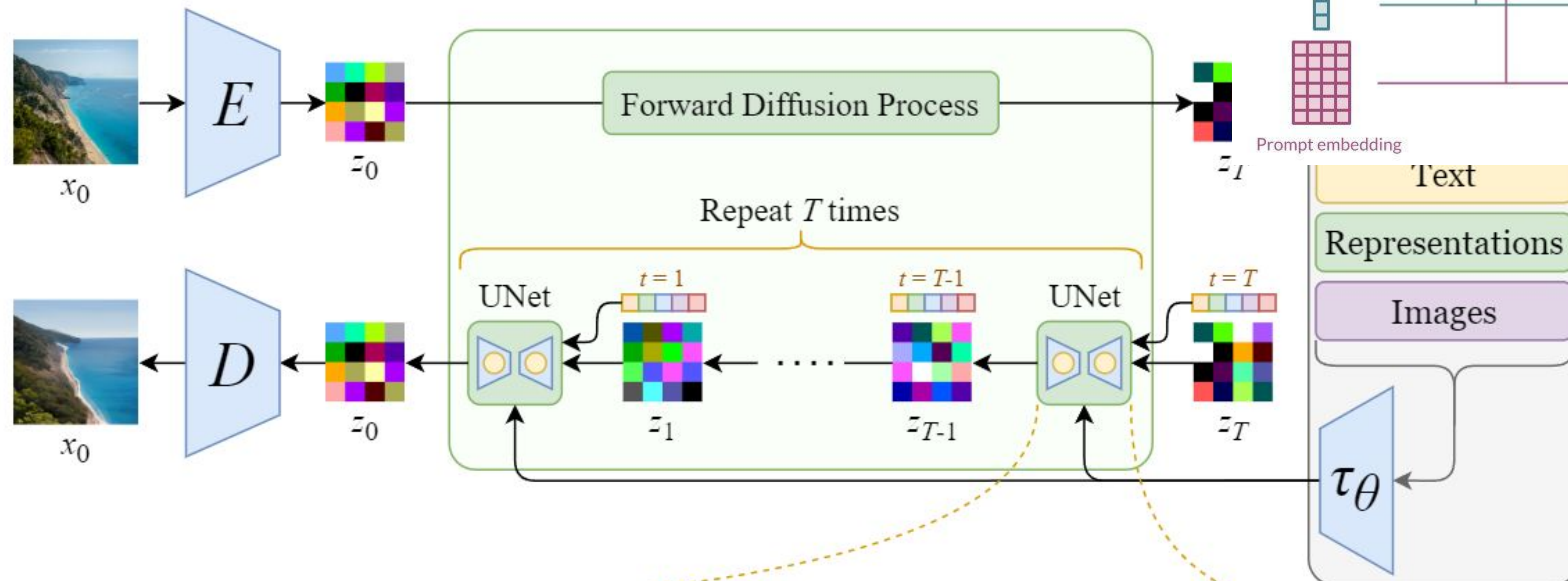
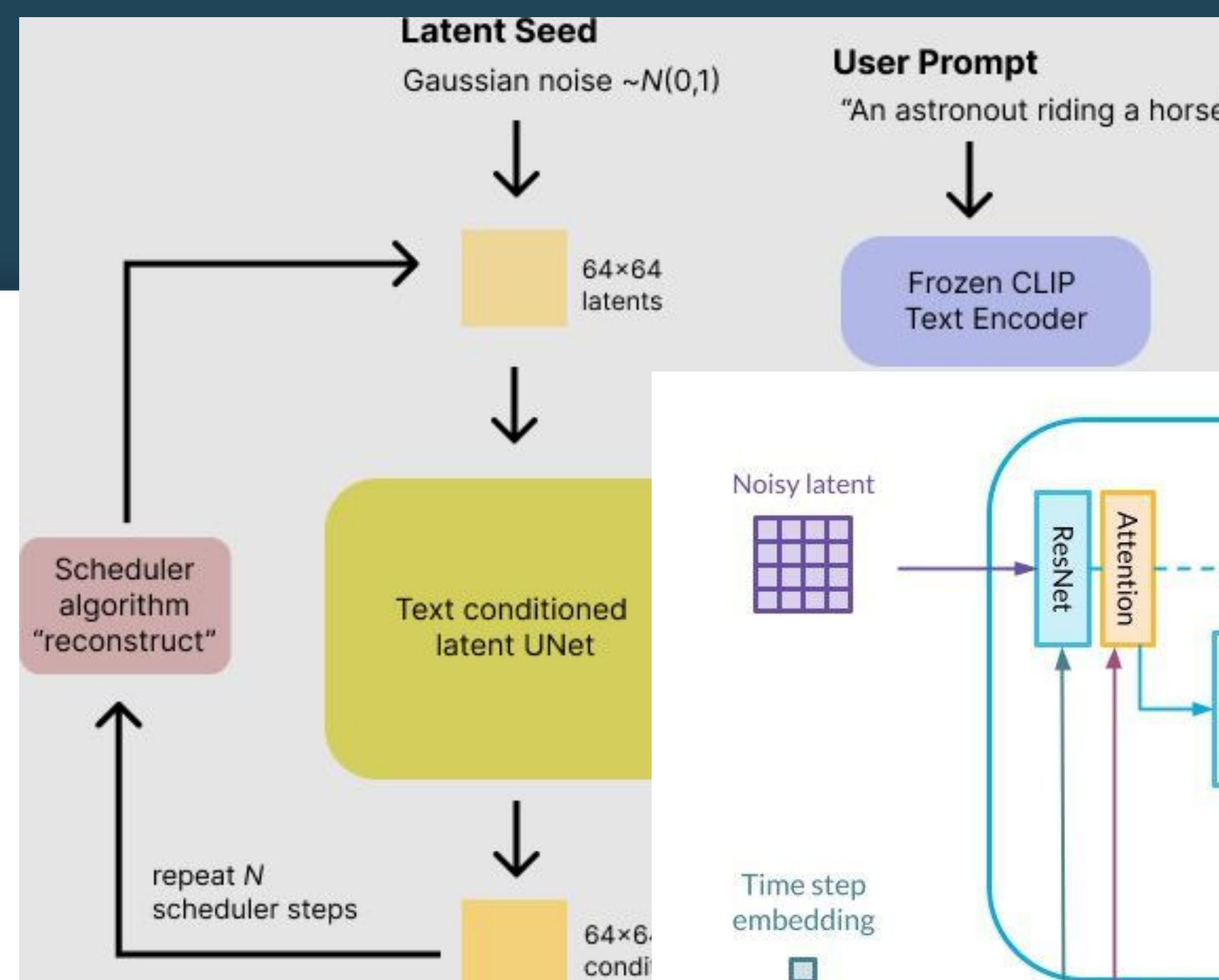
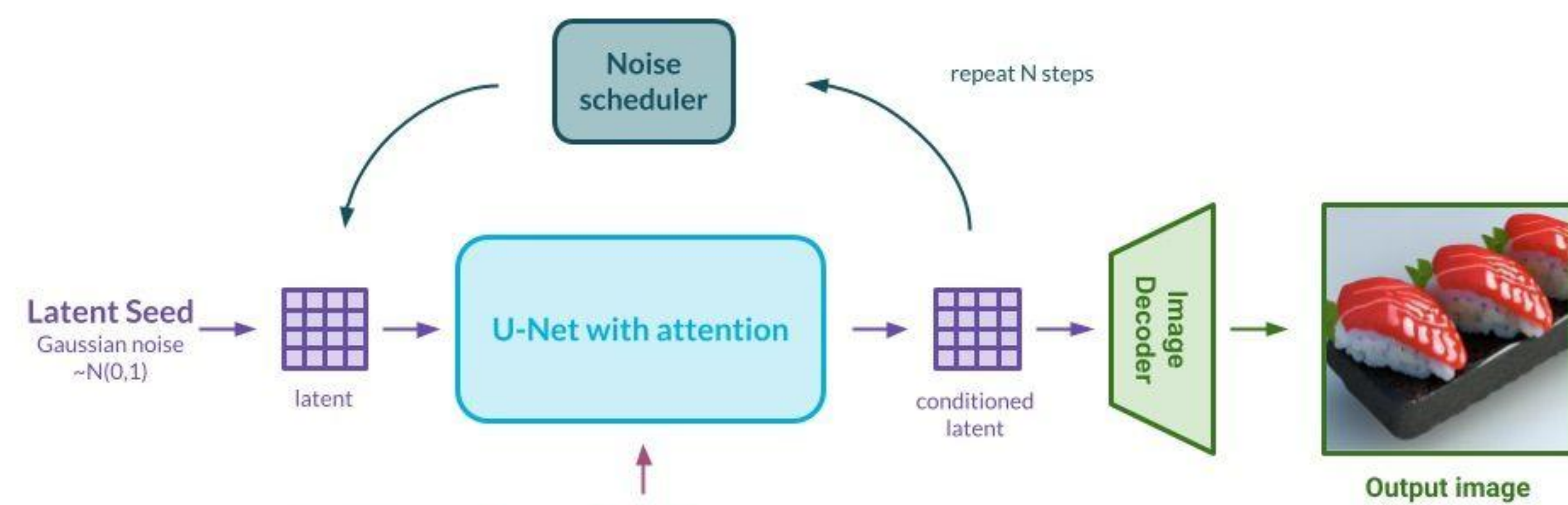


Figure 3. We condition LDMs either via concatenation or by a more general cross-attention mechanism. See Sec. 3.3