
Emergence and Evolution of Interpretable Concepts in Diffusion Models

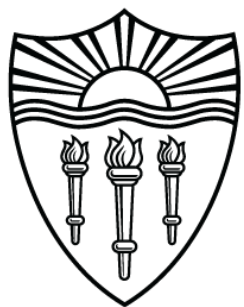
Berk Tinaz* Zalan Fabian* Mahdi Soltanolkotabi

Dept. of Electrical and Computer Engineering

University of Southern California

Los Angeles, CA, USA

tinaz@usc.edu fabian.zalan@gmail.com soltanol@usc.edu



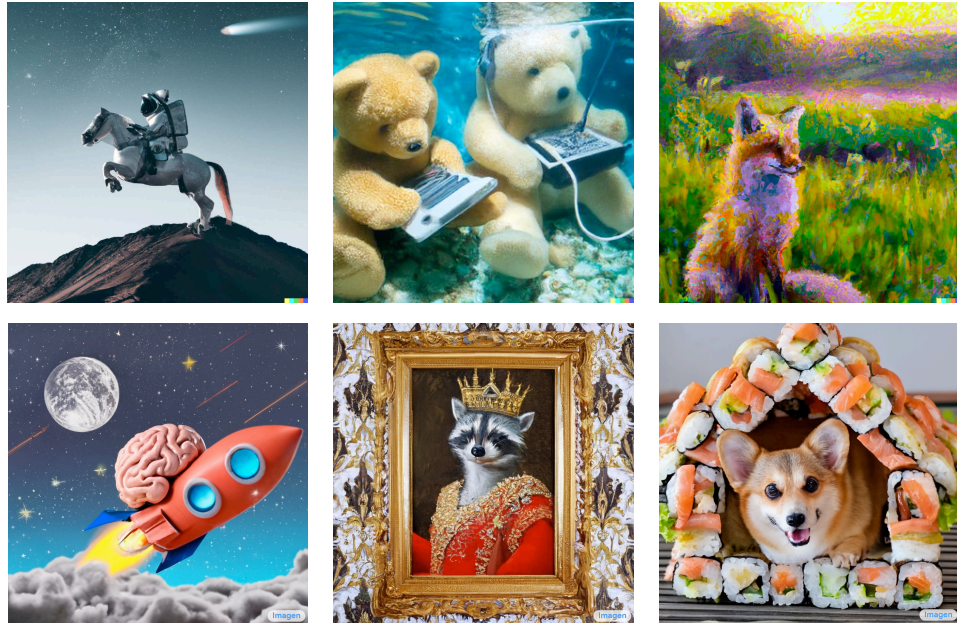
USC

USC Center on AI Foundations for the Sciences

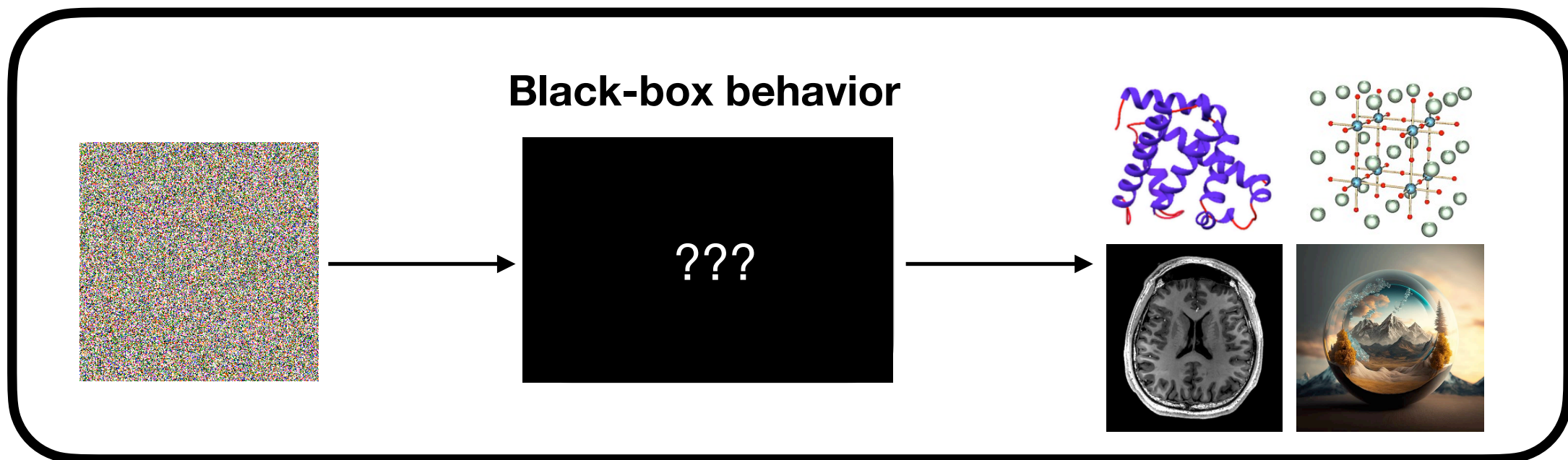


**NEURAL INFORMATION
PROCESSING SYSTEMS**

Interpretable Diffusion



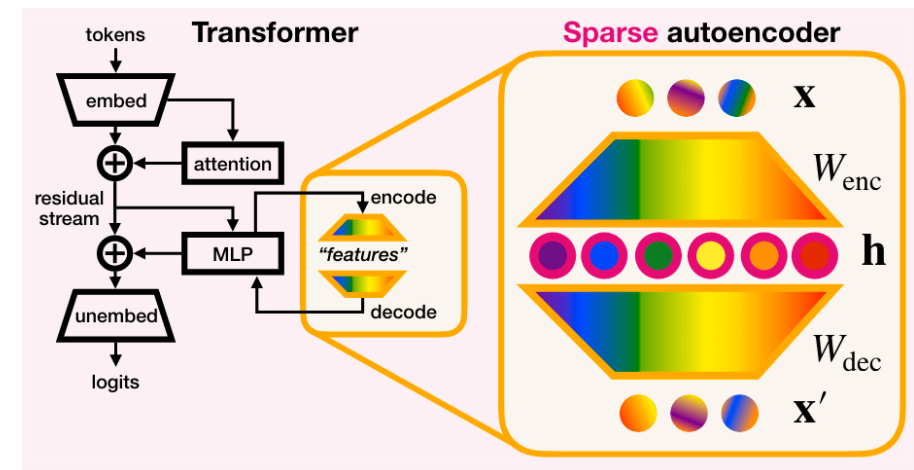
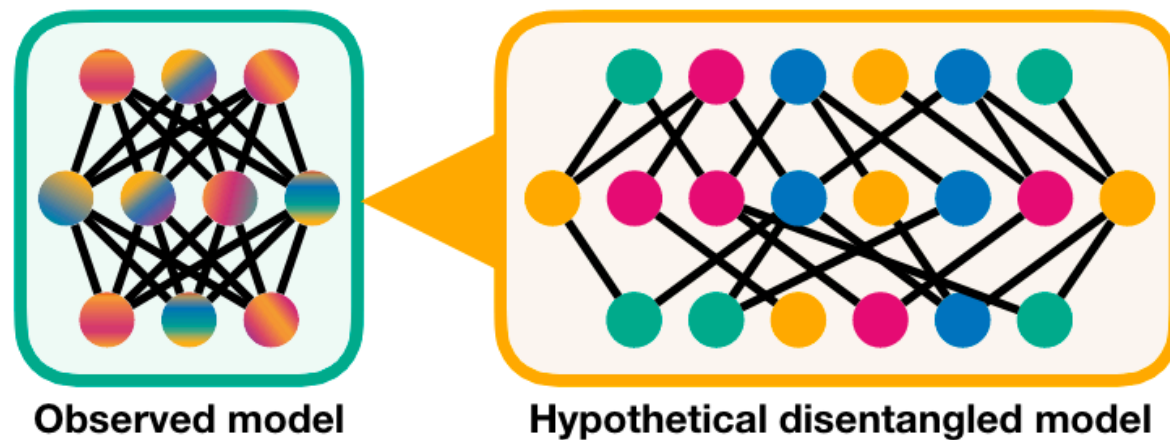
Diffusion models are the **powerhouse** behind modern generative AI.



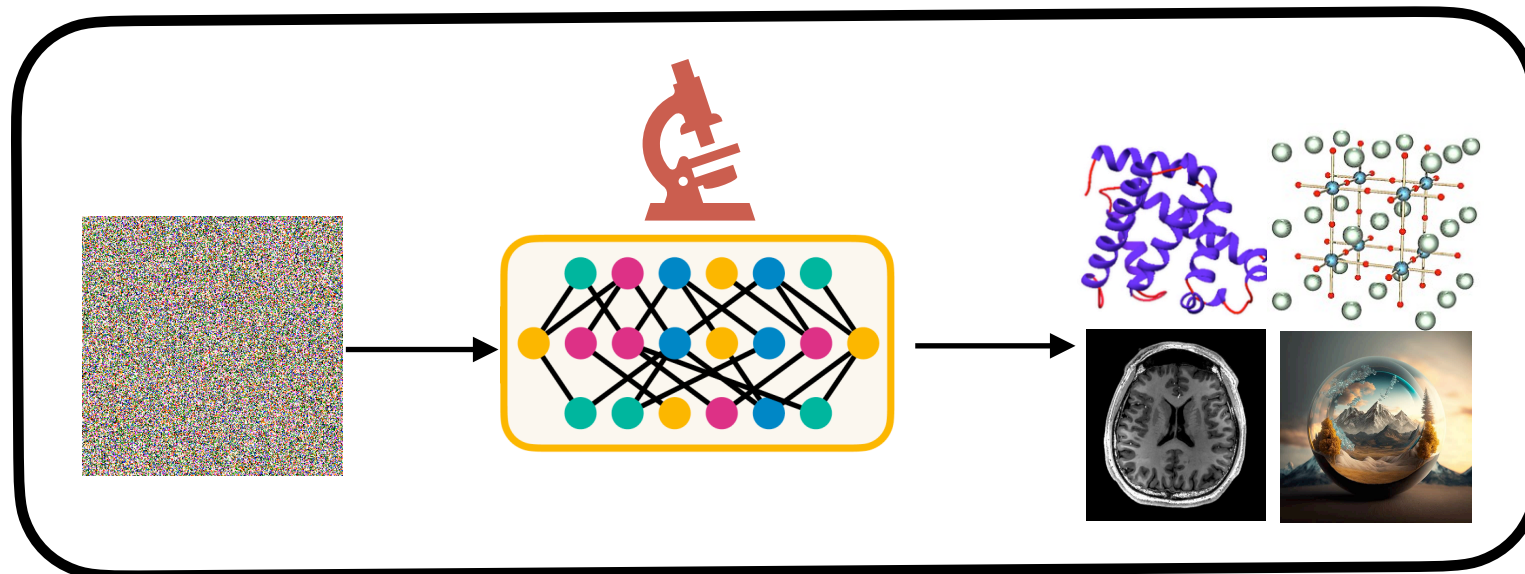
Our understanding of the inner mechanisms of diffusion models remains **limited**.

Sparse Autoencoders

Mechanistic interpretability through Sparse Autoencoders (SAEs)



Diffusion models under the lens of SAEs



Goals:

1. Understand the features diffusion models learn
2. Identify causal mechanisms
3. Achieve fine-grained control of the generative process

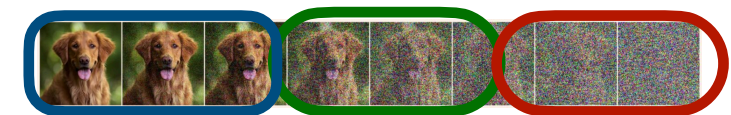
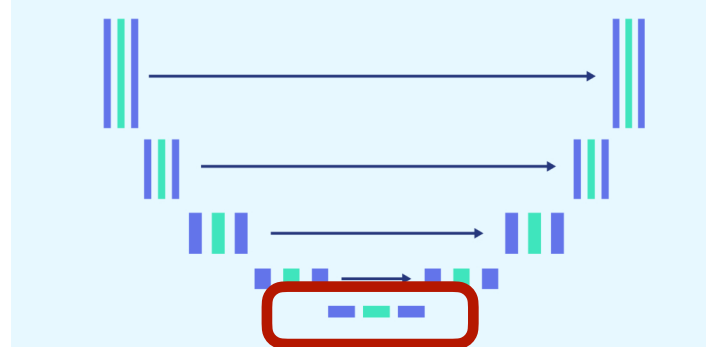
Training Sparse Auto-encoders

Use diverse text prompts and collect activations from the bottleneck layer of **Stable Diffusion v1.4**, at different time steps of the reverse diffusion.

LAION COCO: 600M SYNTHETIC CAPTIONS FROM LAION2B-EN

by: Christoph Schuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, Romain Beaumont, 15 Sep, 2022

Denoising U-Net



Training:

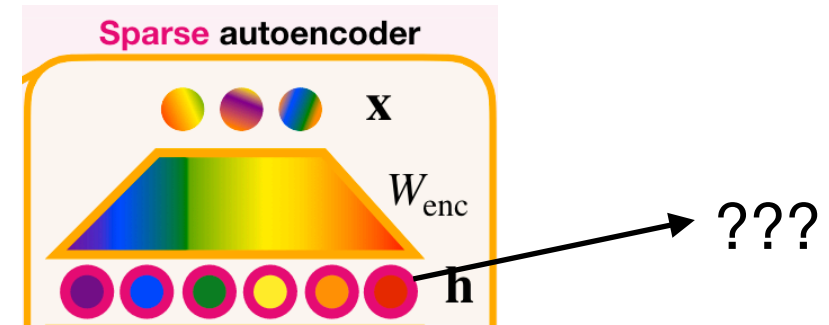
Encoder: $\mathbf{z} = \mathcal{E}_{\theta}(\mathbf{x}) = \text{TopK}(\text{ReLU}(\mathbf{W}_{enc}(\mathbf{x} - \mathbf{b})))$

Decoder: $\hat{\mathbf{x}} = \mathcal{D}_{\theta}(\mathbf{z}) = \mathbf{W}_{dec}\mathbf{z} + \mathbf{b}$

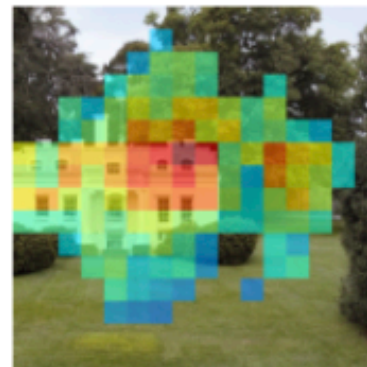
$$\mathcal{L}_{rec}(\theta) = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \alpha \mathcal{L}_{aux}(\theta) \rightarrow \text{Prevent "dead" neurons}$$

Labeling SAE Features

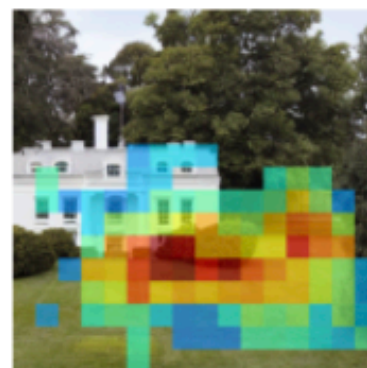
How to automatically label SAE features?



Assign list of objects to each SAE feature using an open-source image segmentation pipeline!



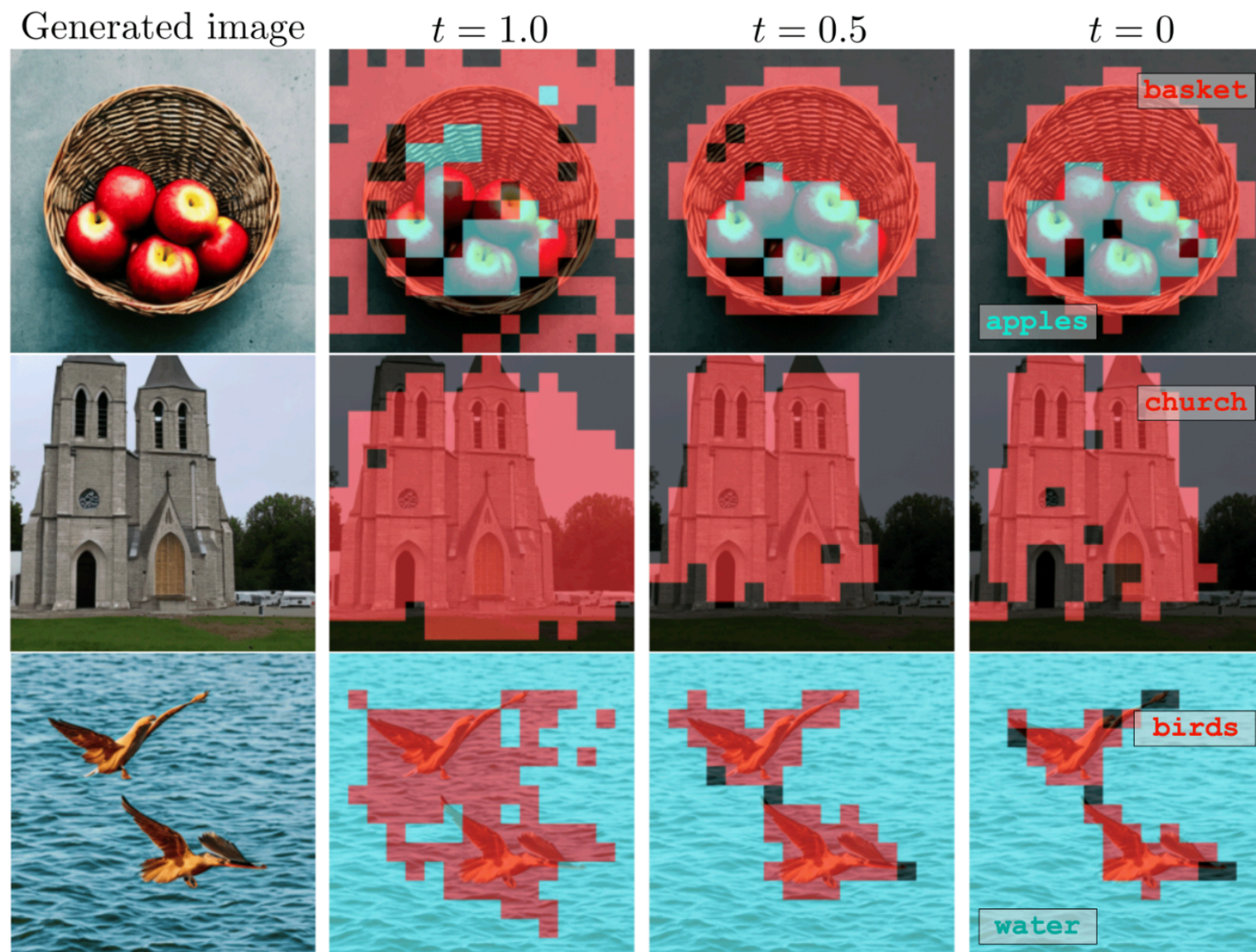
CID: 4657
white house, building



CID: 2577
grass

Semantic Segmentation

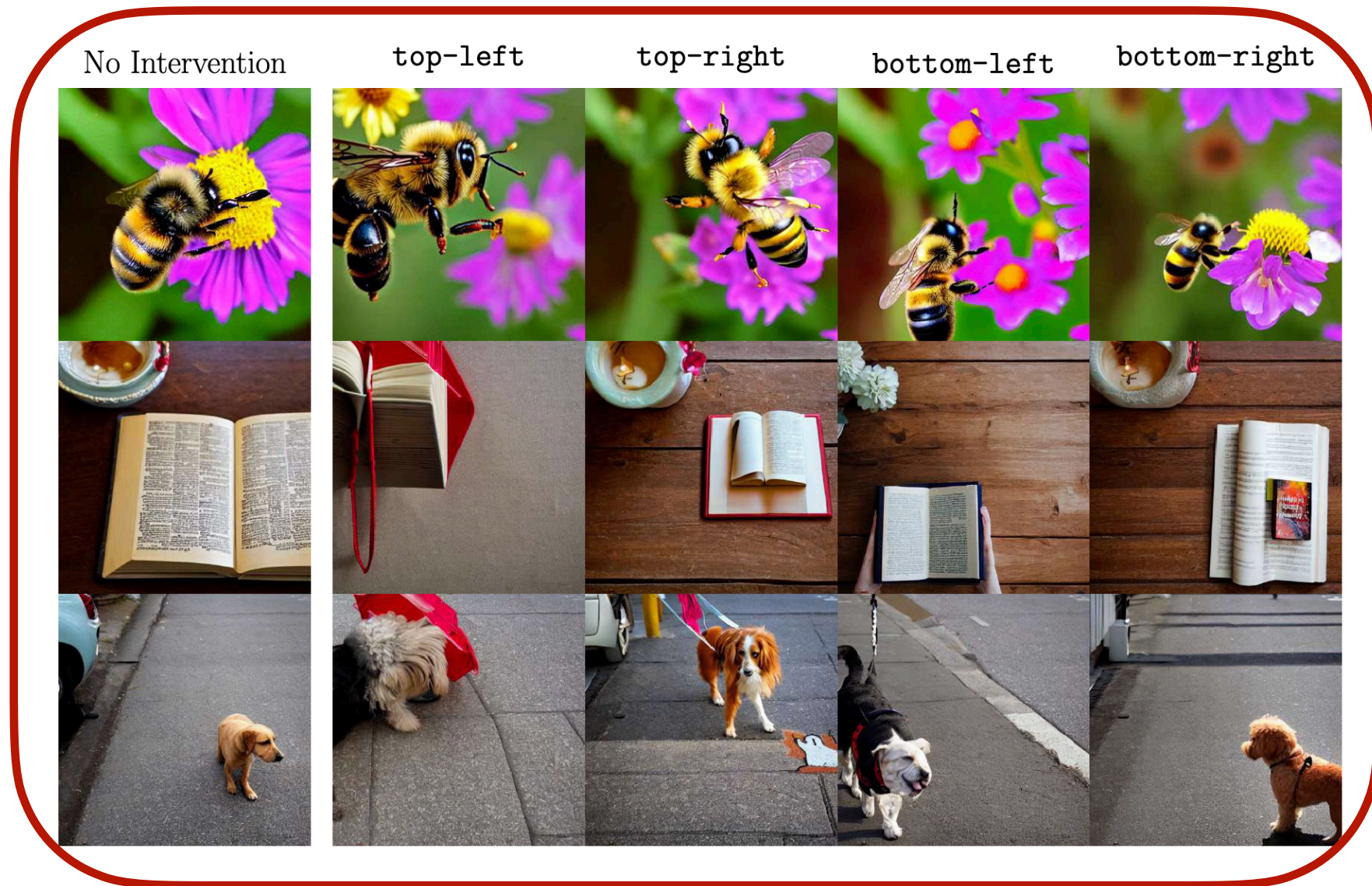
When does the image layout emerge during sampling?



We can **predict** the image layout even **before** the first diffusion step is completed and prediction **gets better** as time progresses.

Spatially Targeted Interventions

Interventions at



We can **restrict** objects to the specified quadrant of the image when intervened in **early stages**.

Spatially Targeted Interventions

Interventions at

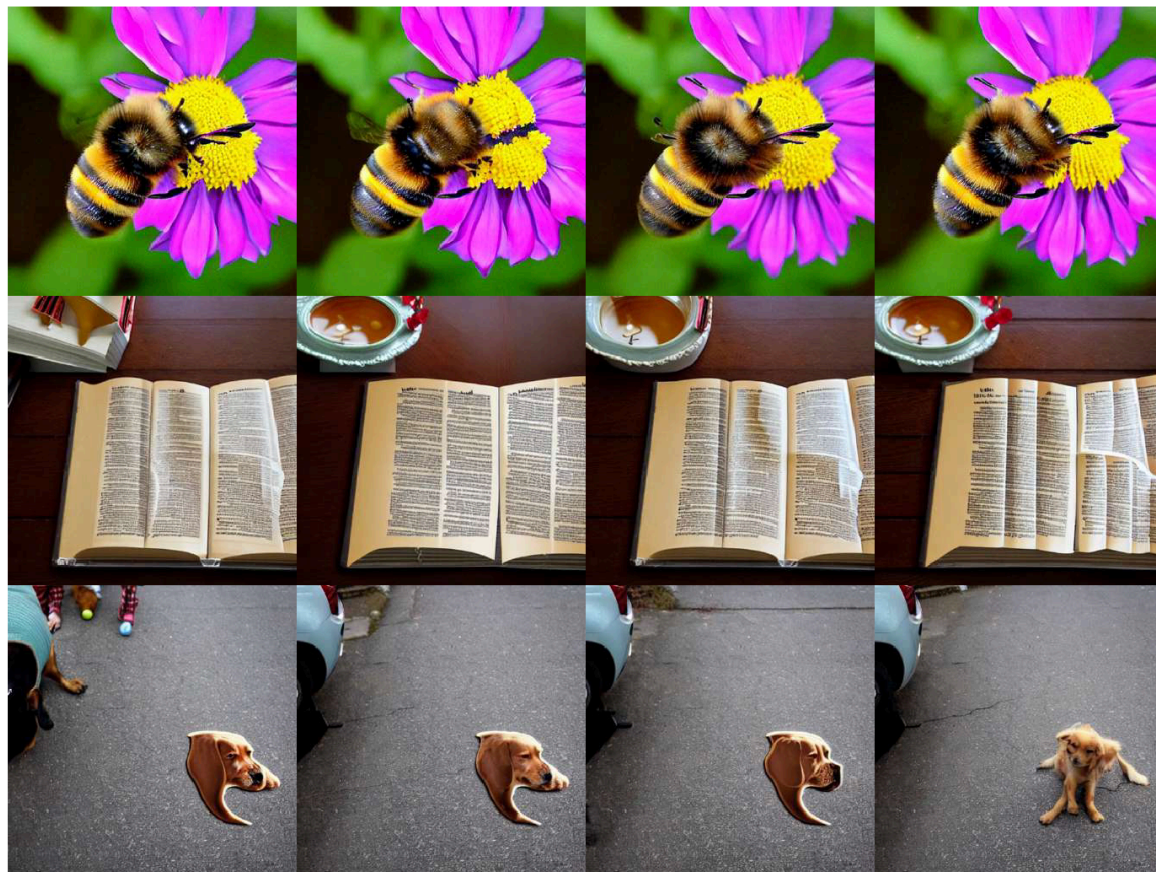


top-left

top-right

bottom-left

bottom-right

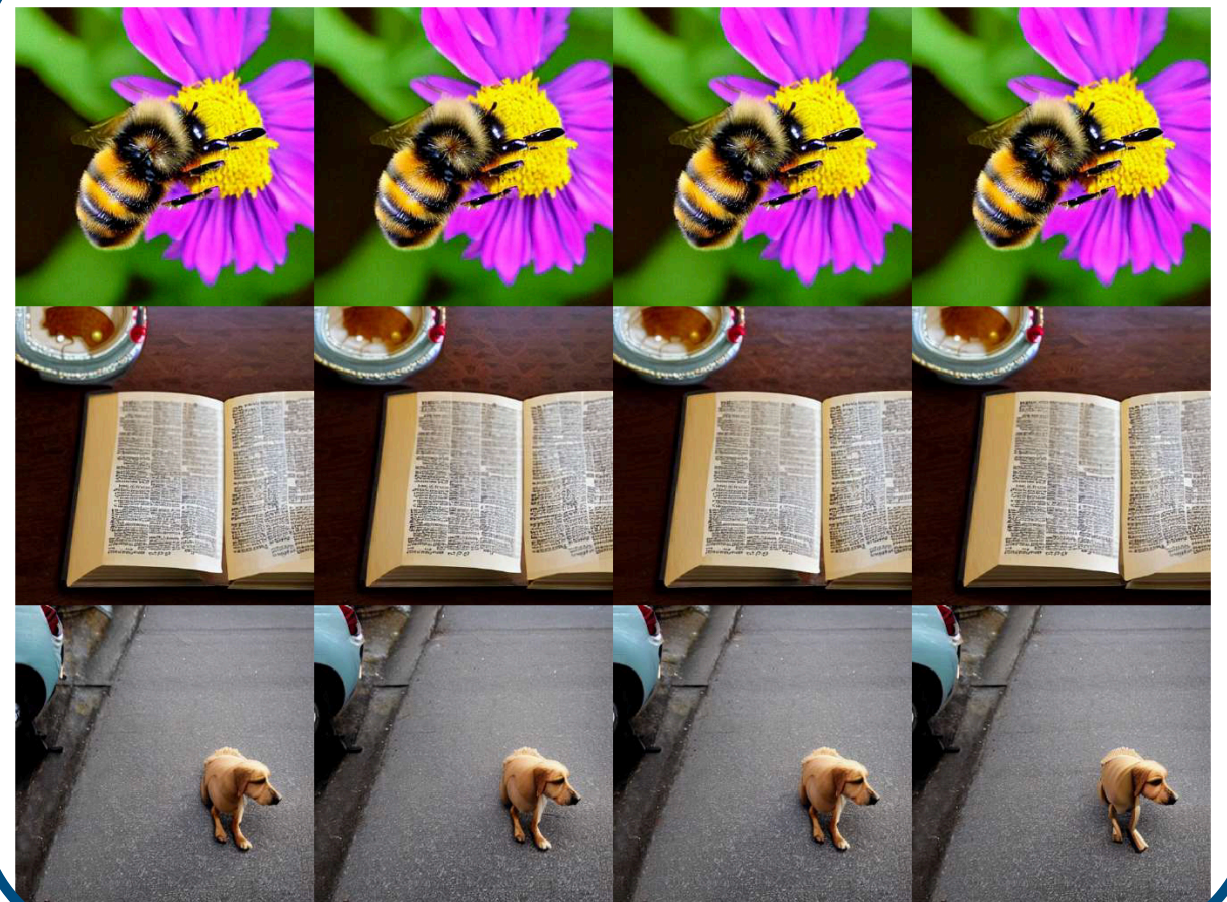


top-left

top-right

bottom-left

bottom-right



Our interventions are **unsuccessful** in the **middle** and **final** stages of reverse diffusion.

Global Interventions (Early Stage)

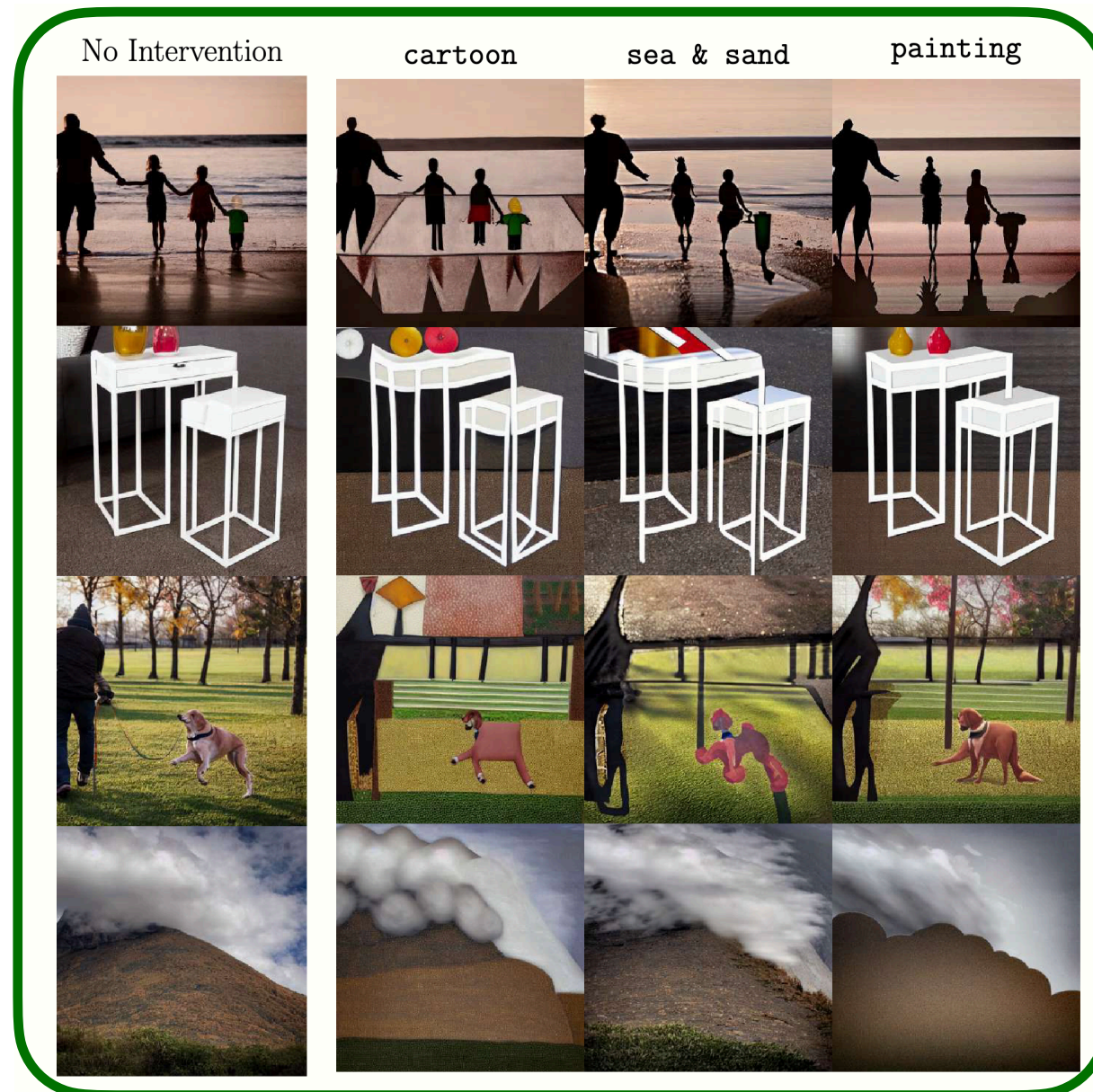
Interventions at



Edits in **initial steps** of reverse diffusion drastically modify the **broad composition** of the image.

Global Interventions (Middle Stage)

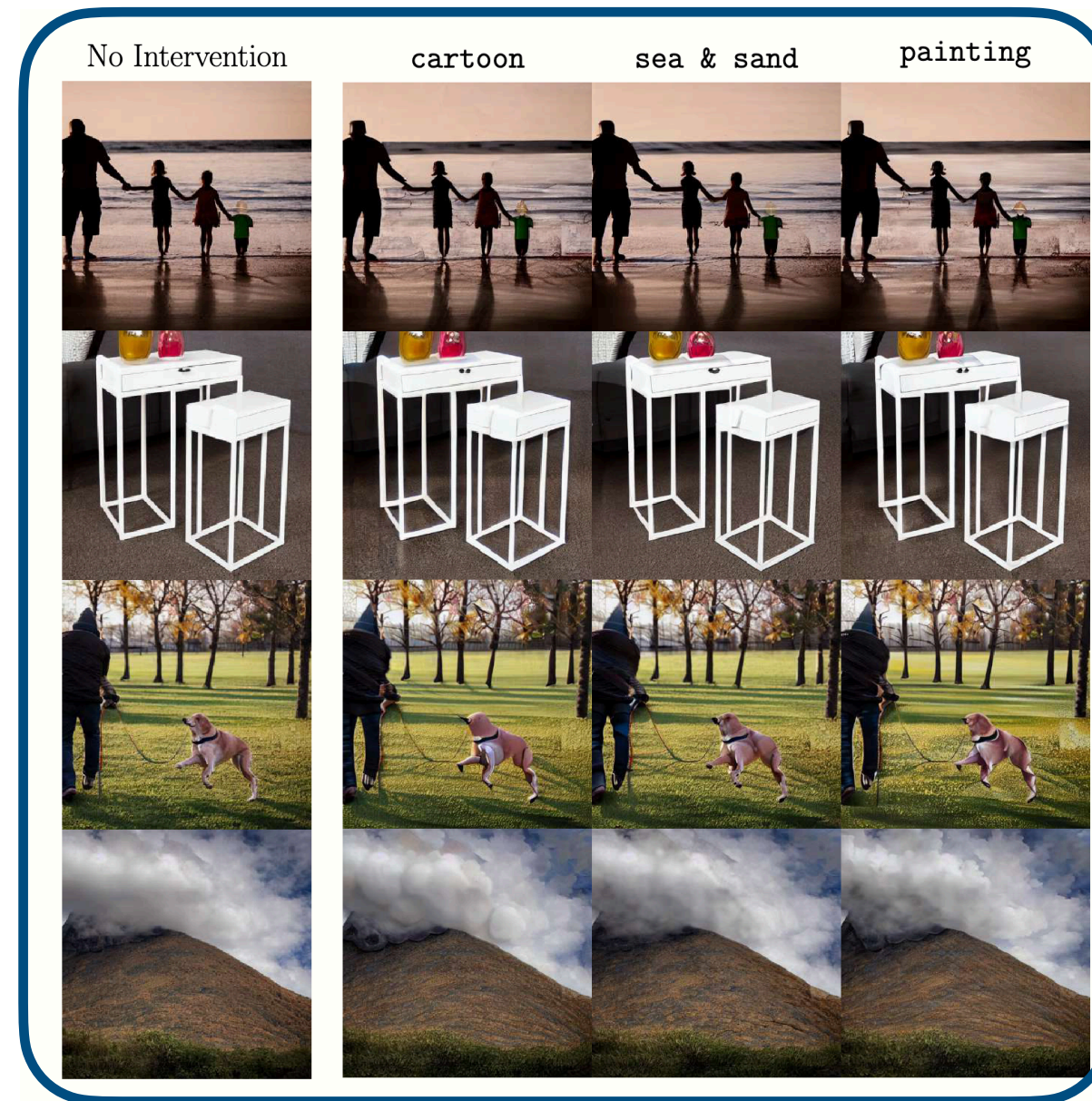
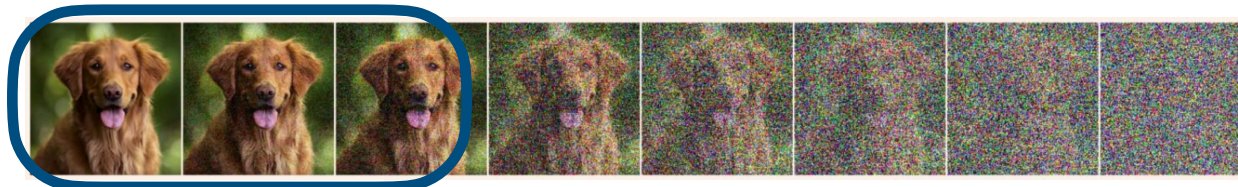
Interventions at



Middle-stage interventions successfully manipulate image **style** without interfering with **image composition**.

Global Interventions (Final Stage)

Interventions at



Global interventions in the **final** stages have **no effect** on style or composition, results only in minor textural changes.

Thank you for your attention!



<https://github.com/berktinaz/stable-concepts>

Paper



Contact: tinaz@usc.edu