



LaRender: Training-Free Occlusion Control in Image Generation via Latent Rendering

ICCV 2025 (*Oral & Award Candidate*)



Xiaohang Zhan



Dingming Liu

Tencent

Speaker: Xiaohang Zhan
(Tencent => Adobe)

Text-to-image models struggle in occlusion control

SDXL



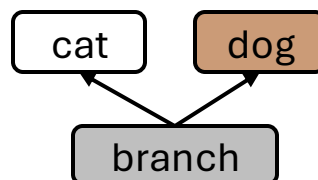
FLUX



Nano Banana

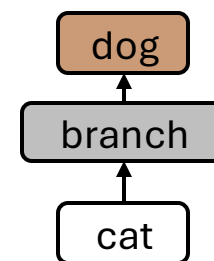
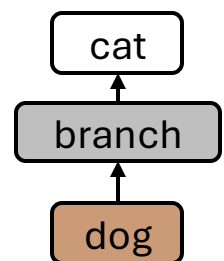


Prompt: In a forest, a white cat and a brown dog are standing, a long branch is in front of them and occludes both of them.



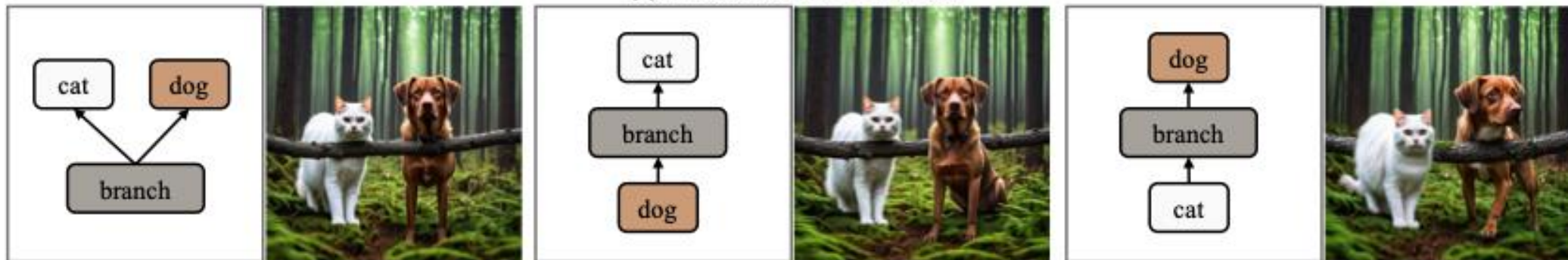
Text-to-image models struggle in occlusion control

Nano Banana

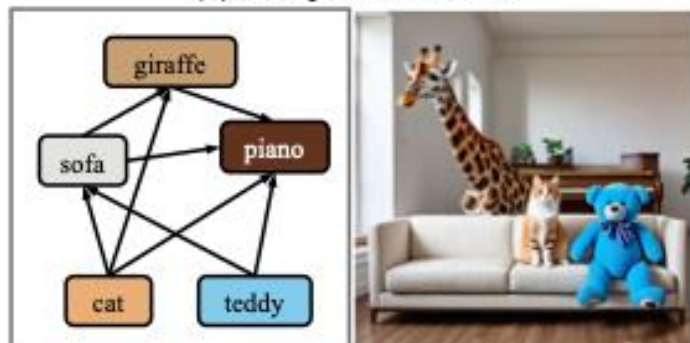


What does LaRender do?

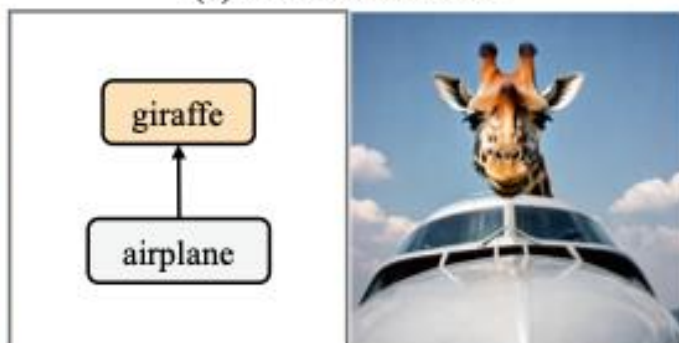
(1) Precise control of occlusion



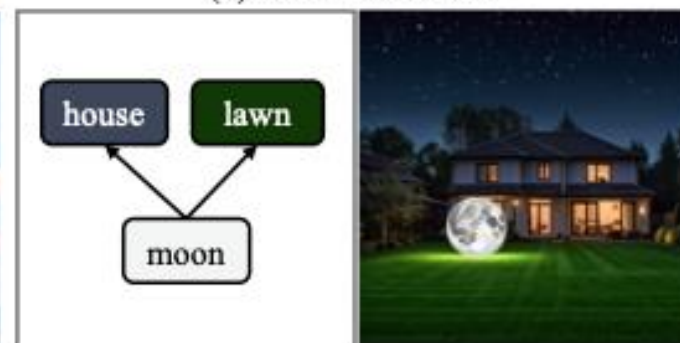
(2) Complex occlusion



(3) Unusual occlusion



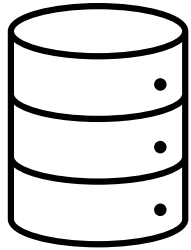
(4) Surreal occlusion



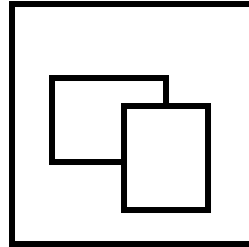
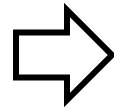
(5) Change opacity



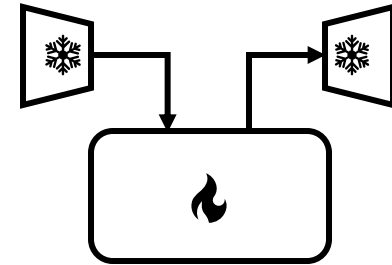
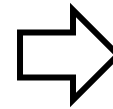
A popular way



Data curation



Control signal design

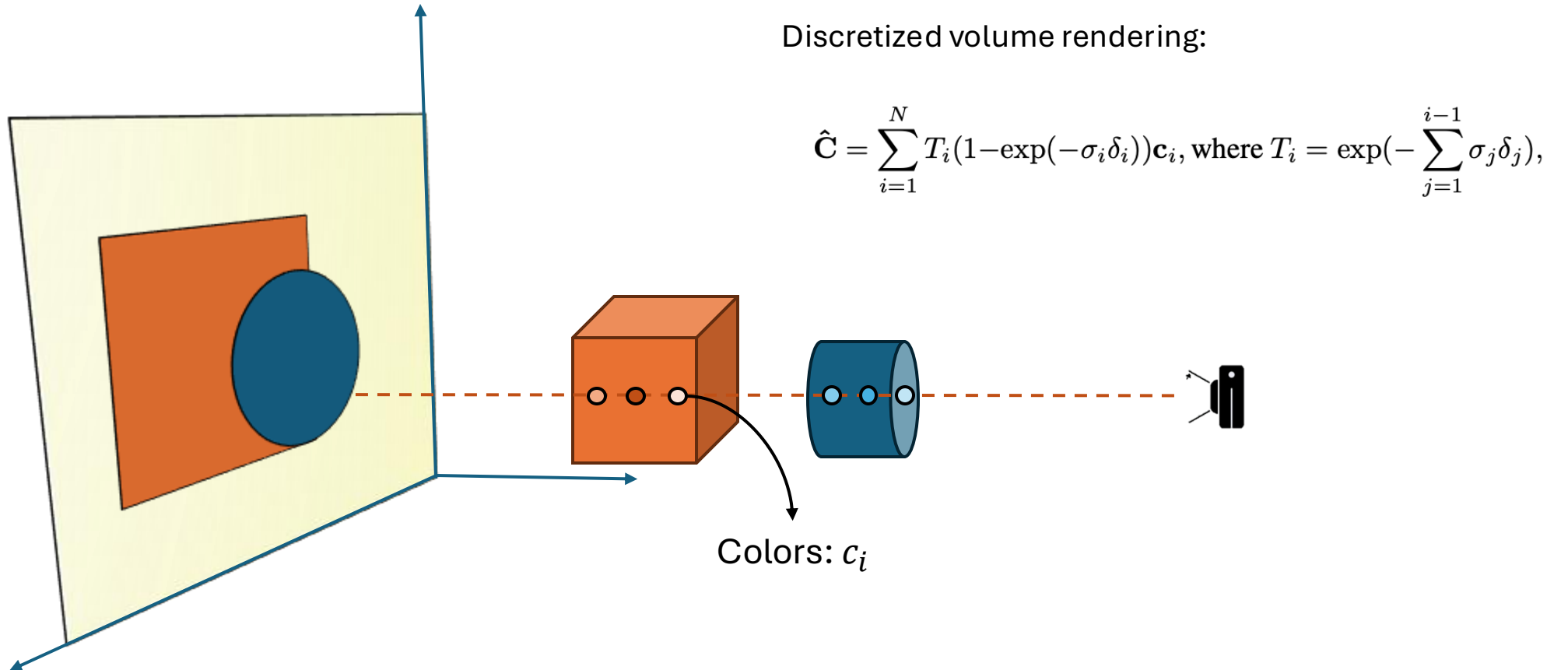


Model finetuning



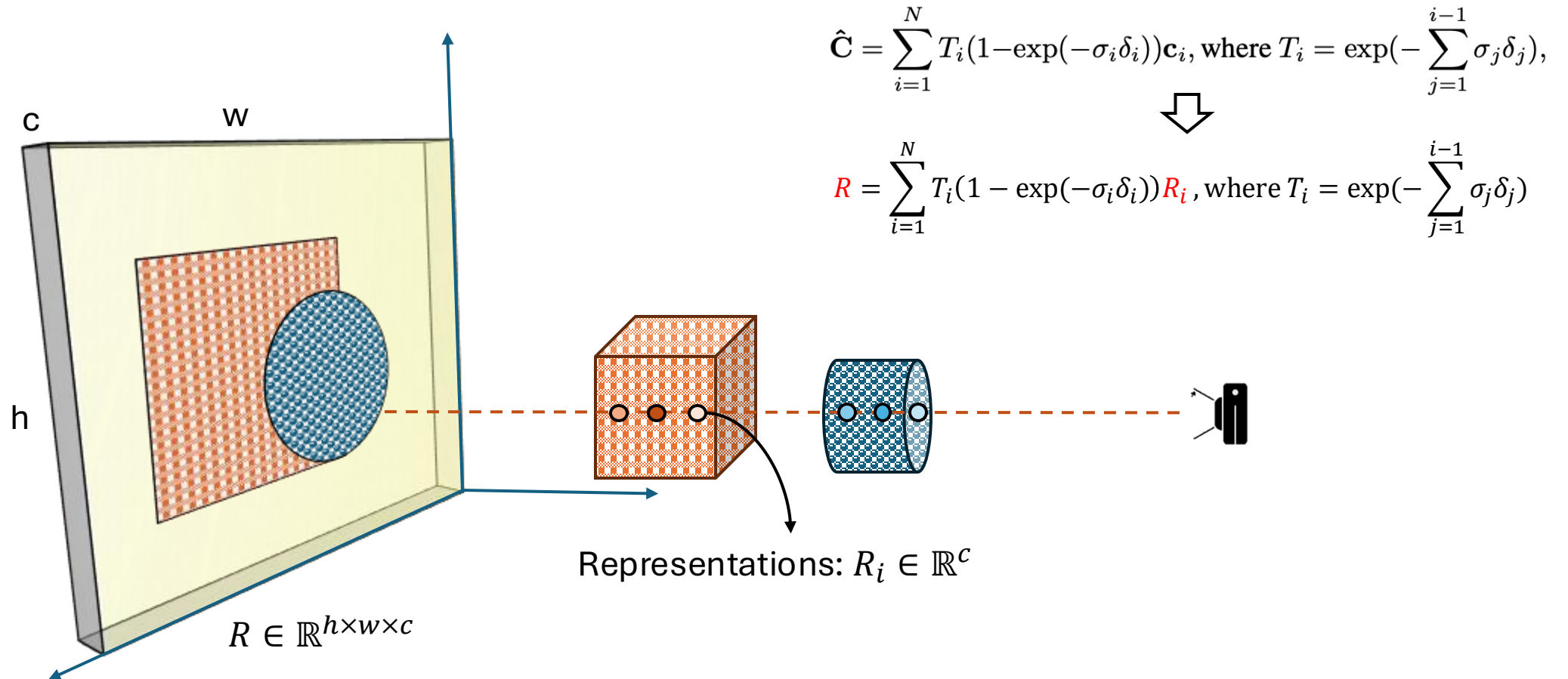
It would work, but the cost is high.

What is visual occlusion?



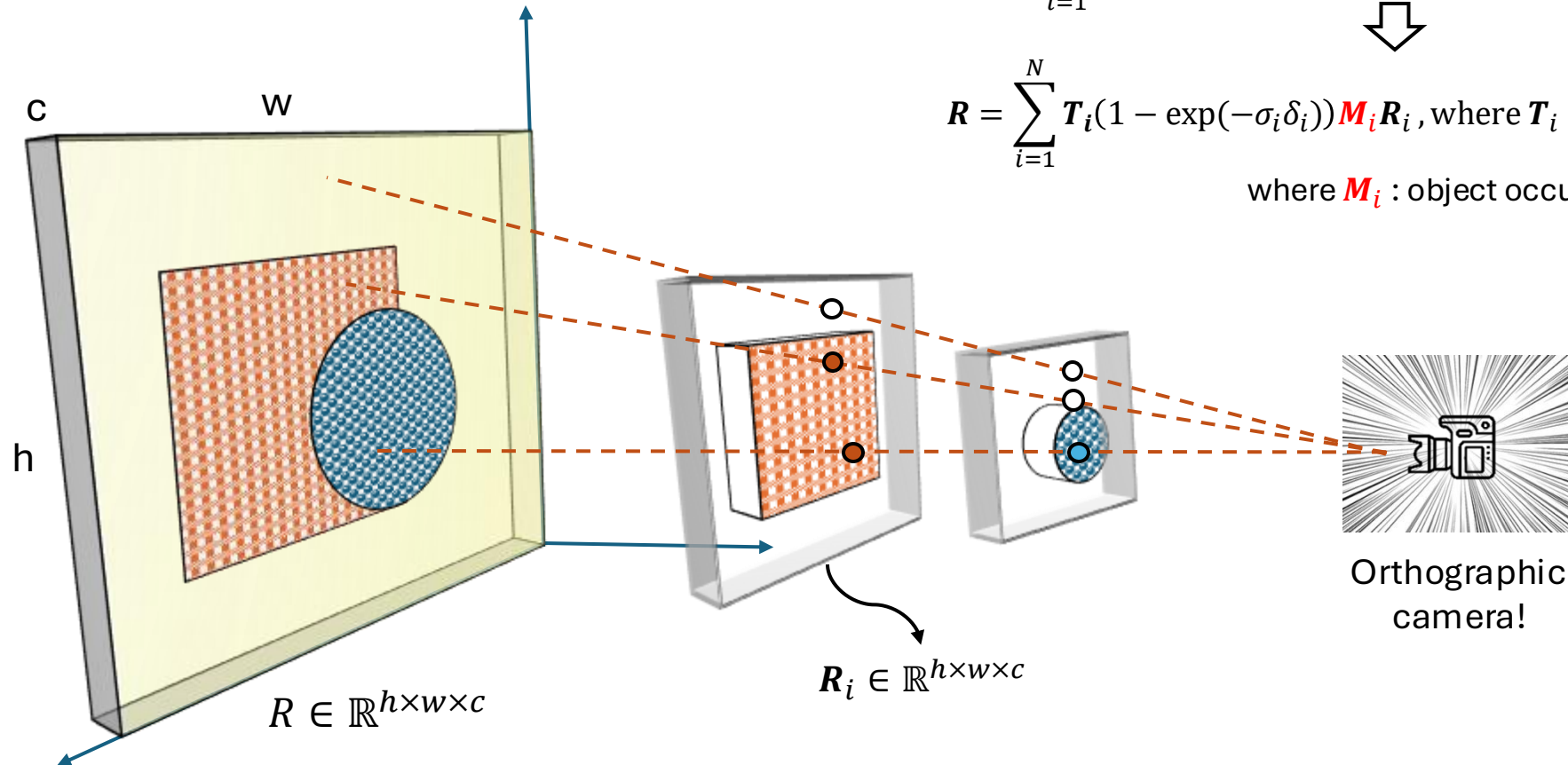
What if they are high dimensional representations?

From physical rendering to latent rendering



3D representations are more difficult to obtain than 2D ones. Let's simplify them.

From physical rendering to latent rendering



$$R = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) R_i, \text{ where } T_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j)$$



$$R = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{M}_i R_i, \text{ where } T_i = \exp(-\sum_{j=1}^{i-1} \mathbf{M}_j \sigma_j \delta_j)$$

where \mathbf{M}_i : object occupancy mask

Using 2D hidden states from text-to-image diffusion models while keeping 3D layouts and the shape of the cross sections.



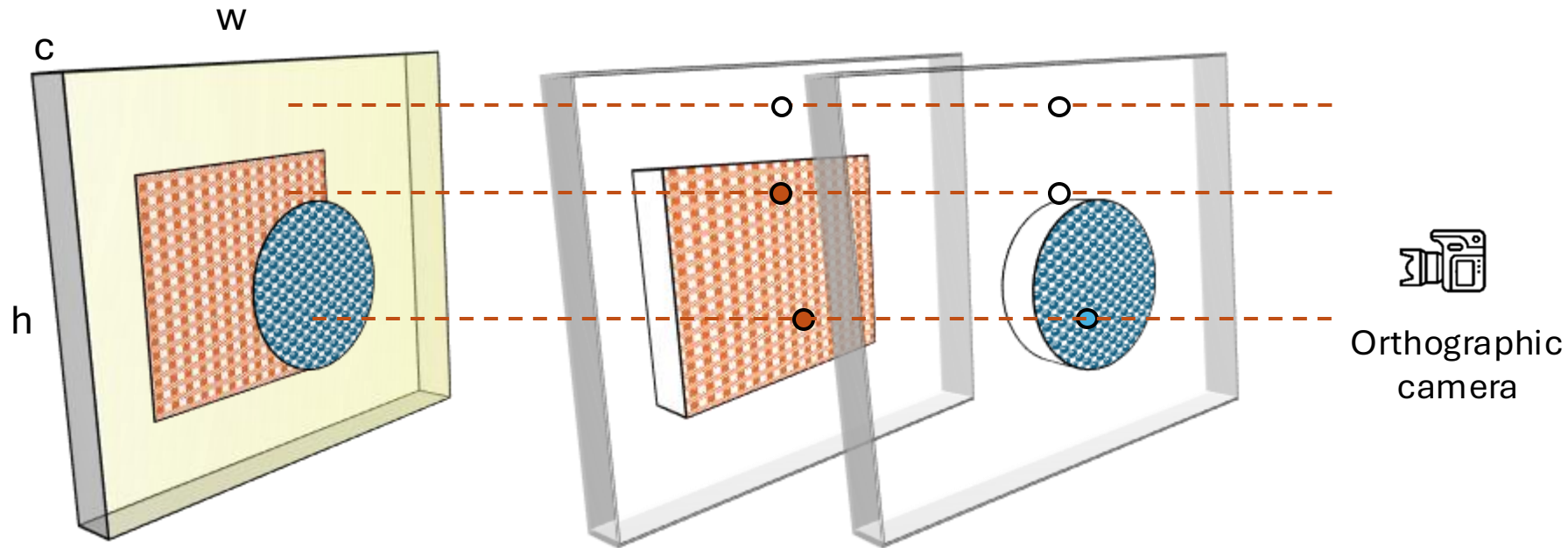
But how to determine their distances to the camera, and the FOV of the camera?

From physical rendering to latent rendering

$$R = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) M_i R_i, \text{ where } T_i = \exp\left(-\sum_{j=1}^{i-1} M_j \sigma_j \delta_j\right)$$

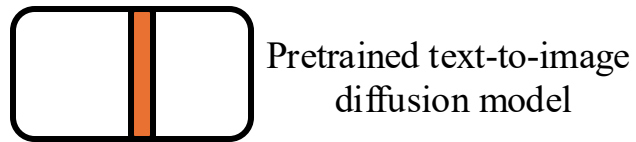


$$R = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i)) M_i R_i, \text{ where } T_i = \exp\left(-\sum_{j=1}^{i-1} M_j \sigma_j\right)$$

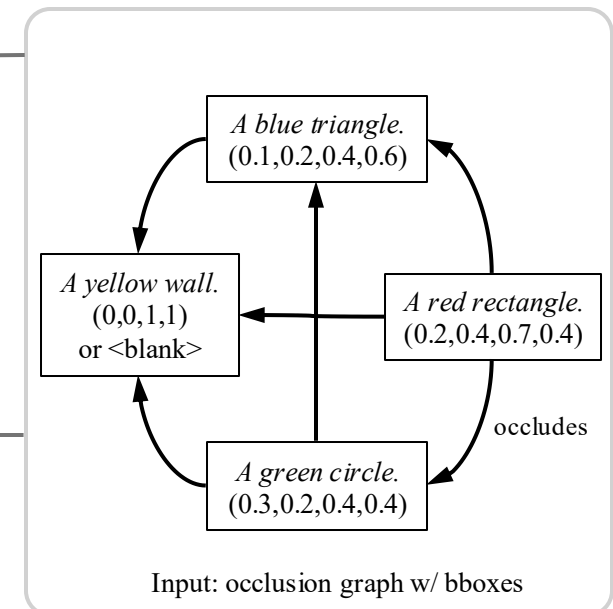
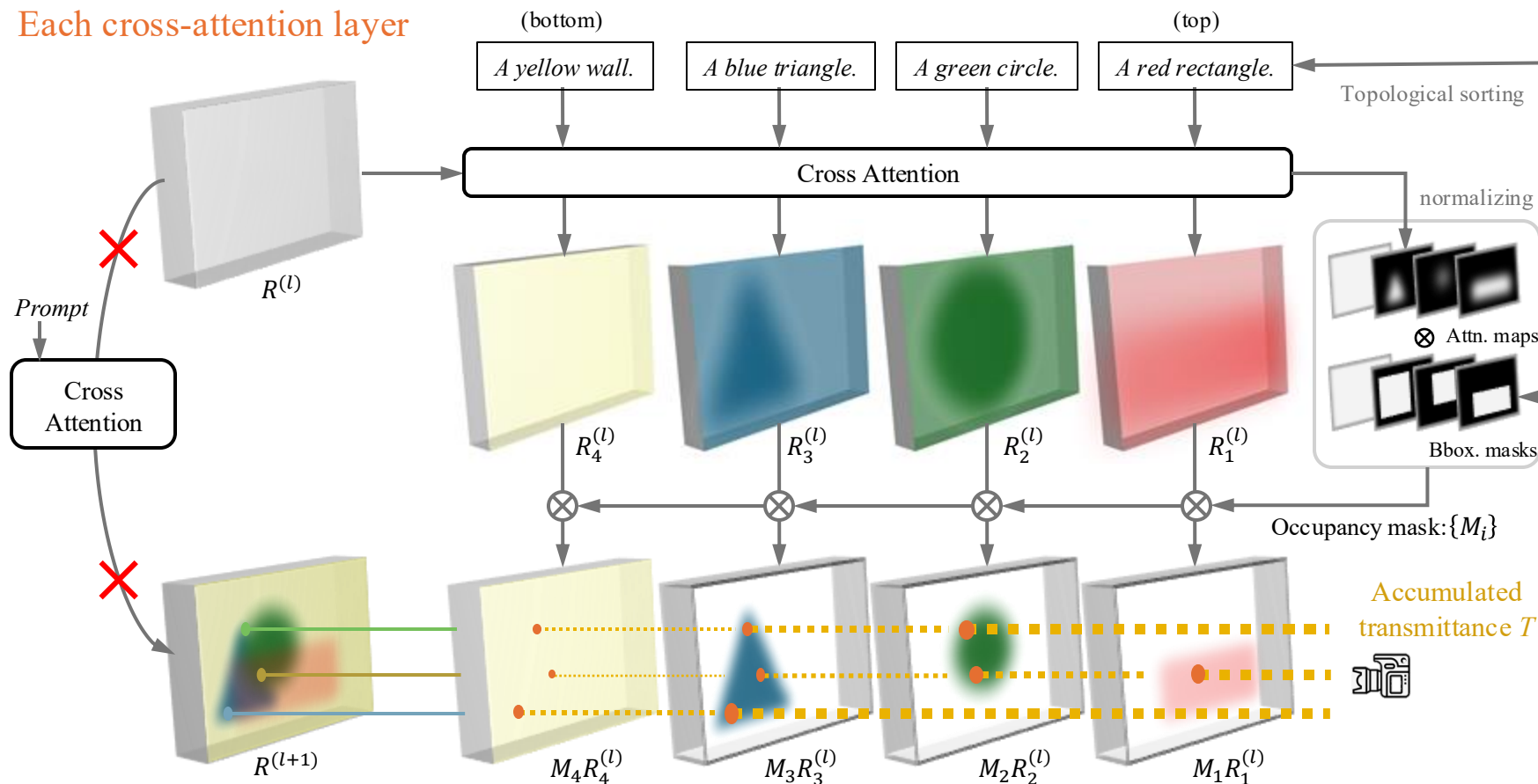


Now we created the foundation of Latent Rendering!

LaRender: a non-parametric framework



Each cross-attention layer



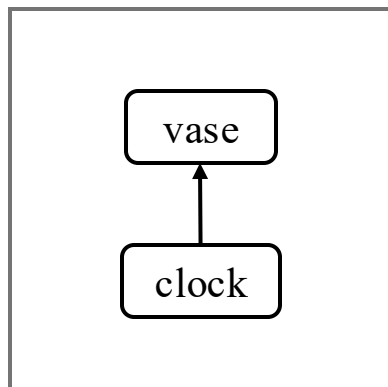
$$\mathbf{R}^{(l+1)} = \frac{1}{\mathbf{S}} \sum_{i=1}^N \mathbf{T}_i (1 - \exp(-\sigma_i)) \mathbf{M}_i \mathbf{R}_i^{(l)},$$

$$\mathbf{S} = \sum_{i=1}^N \mathbf{T}_i (1 - \exp(-\sigma_i)) \mathbf{M}_i,$$

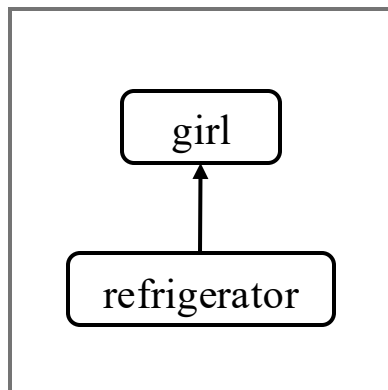
$$\mathbf{T}_i = \exp \left(- \sum_{j=1}^{i-1} \mathbf{M}_j \sigma_j \right), \quad \sigma_i \in [0, +\infty)$$

LaRender results

Prompt: A vase hidden by a clock.



Prompt: A girl hidden by a refrigerator.



SDXL
(text control)

FLUX.1-dev
(text control)

MIGC
(layout control)

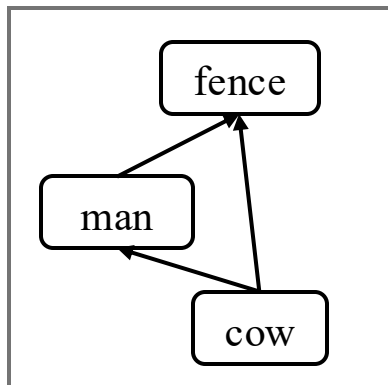
3DIS
(layout control)

LaRender
(occlusion control)

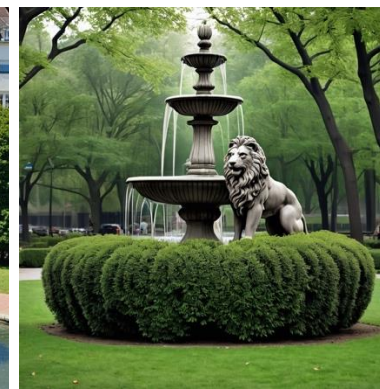
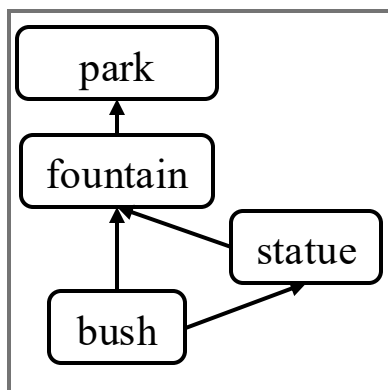
Comparison

LaRender results

Prompt: A cow occludes a man and a fence, the man occludes the fence.



Prompt: In a park, a fountain is partially obscured by a lion statue, and in front of them is a bush that hides both of them.



SDXL
(text control)

FLUX.1-dev
(text control)

MIGC
(layout control)

3DIS
(layout control)

LaRender
(occlusion control)

Comparison

LaRender results

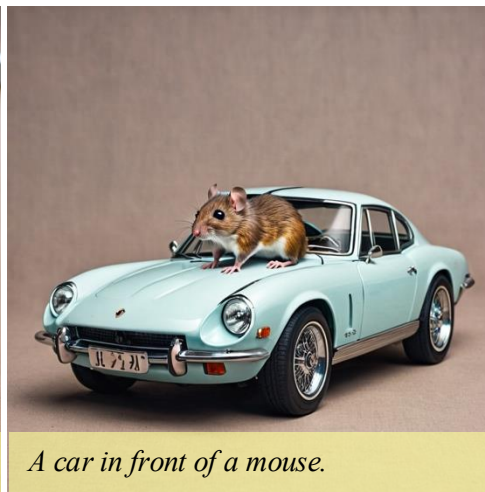


Similar layout, different occlusion

LaRender results



Inaccurate position



Wrong occlusion



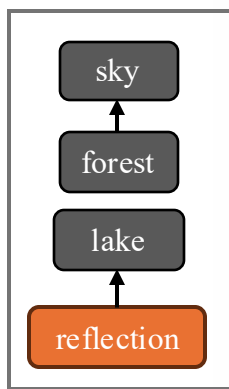
Concept lost



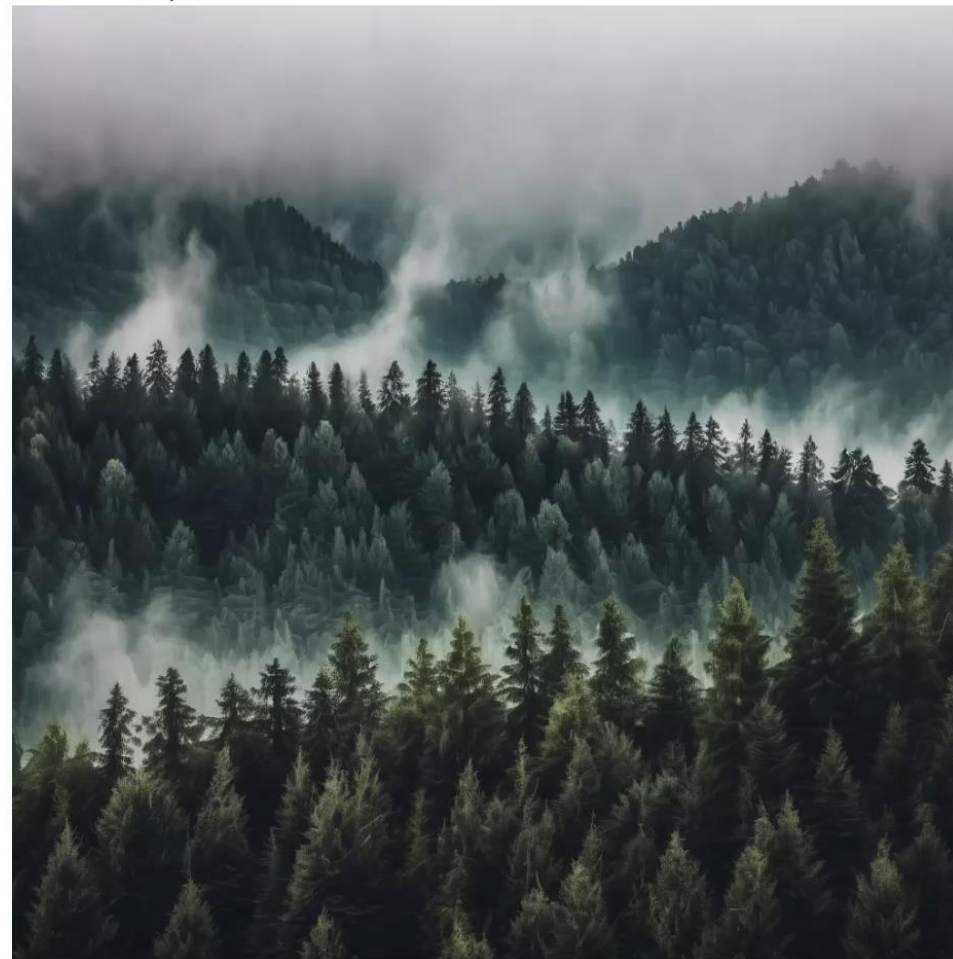
Concept mixed

Failure cases

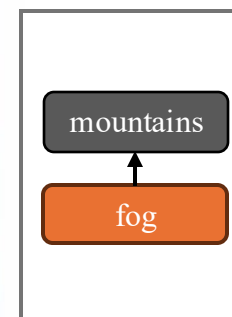
Semantic opacity control



reflection: 0.05

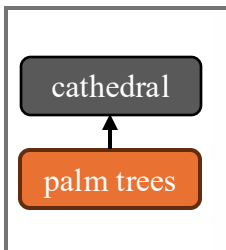


fog: 0.10



$$\alpha_i = 1 - \exp(-\sigma_i) \in [0, 1)$$

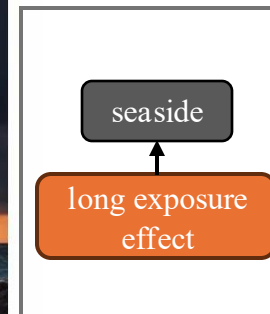
Semantic opacity control



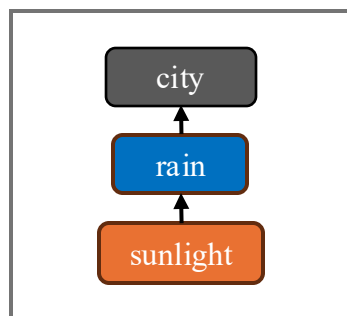
palm trees: 0.00



long exposure effects: 0.00



Semantic opacity control

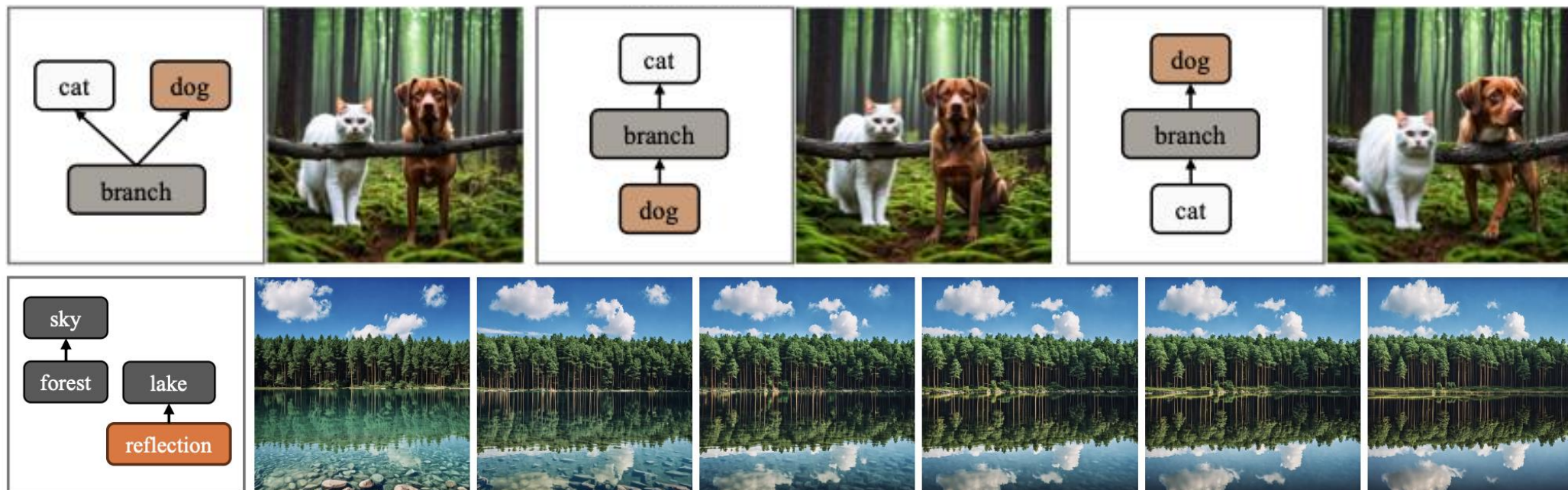


Rain: 0.00

Sunlight: 0.00



Q&A



LaRender provides training-free control over object occlusions and effect opacity in image generation.



Project page (code available)