

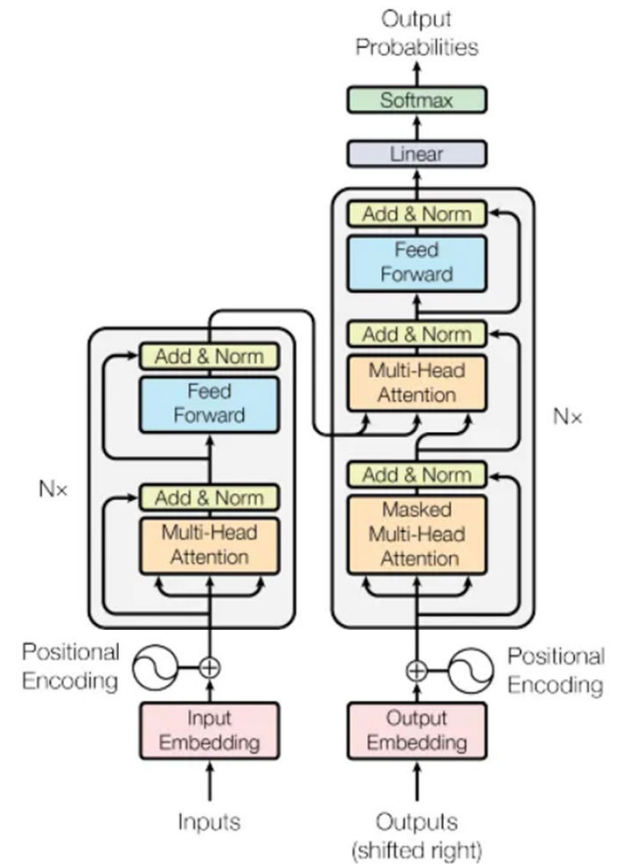
Attention Please

Attention Mechanism in Generative Models

Daniel Cohen-Or

“Attention is all you need” , Vaswani et al. 2017

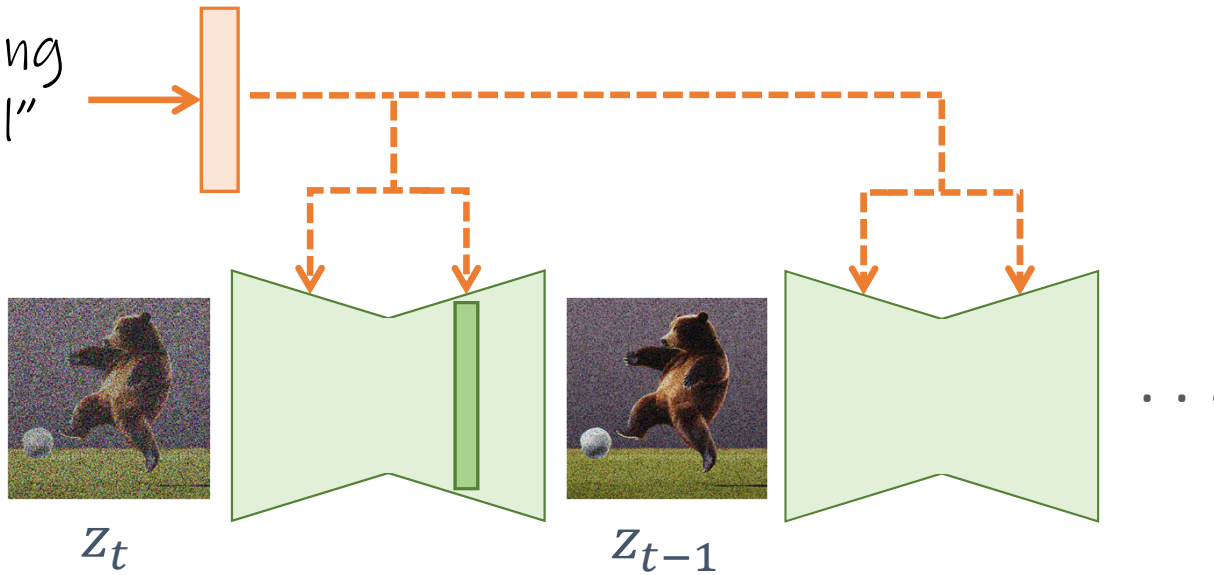
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



From “Attention is all you need” paper by Vaswani, et al., 2017 [1]

The Denoising Process with Attention Layers

"A bear kicking a soccer ball"



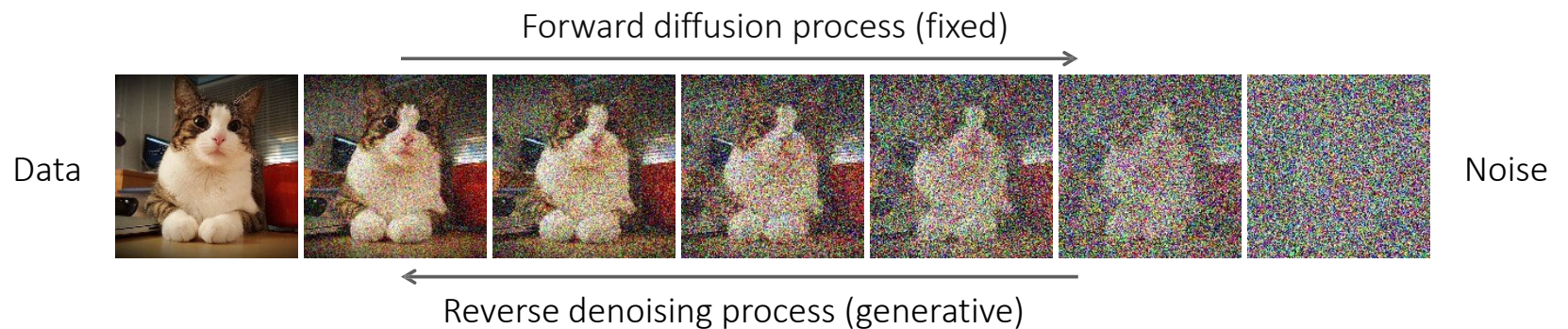
Denoising network

Denoising Diffusion Models

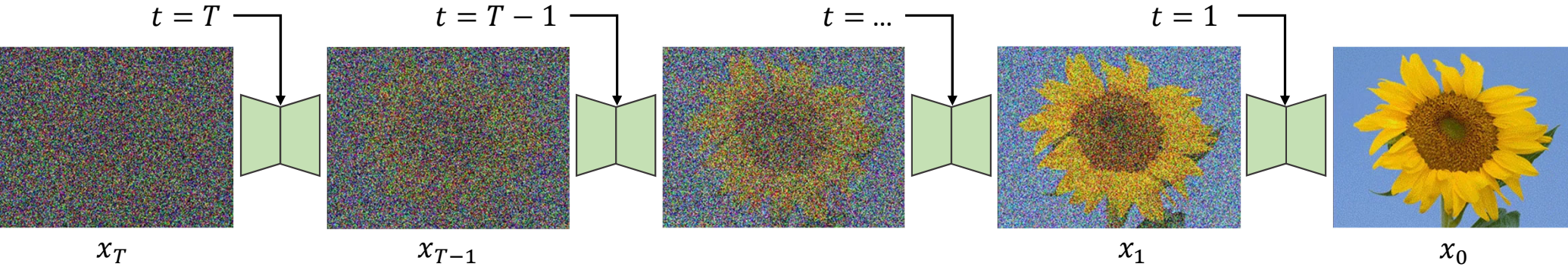
Learning to generate by denoising

Denoising diffusion models consist of two processes:

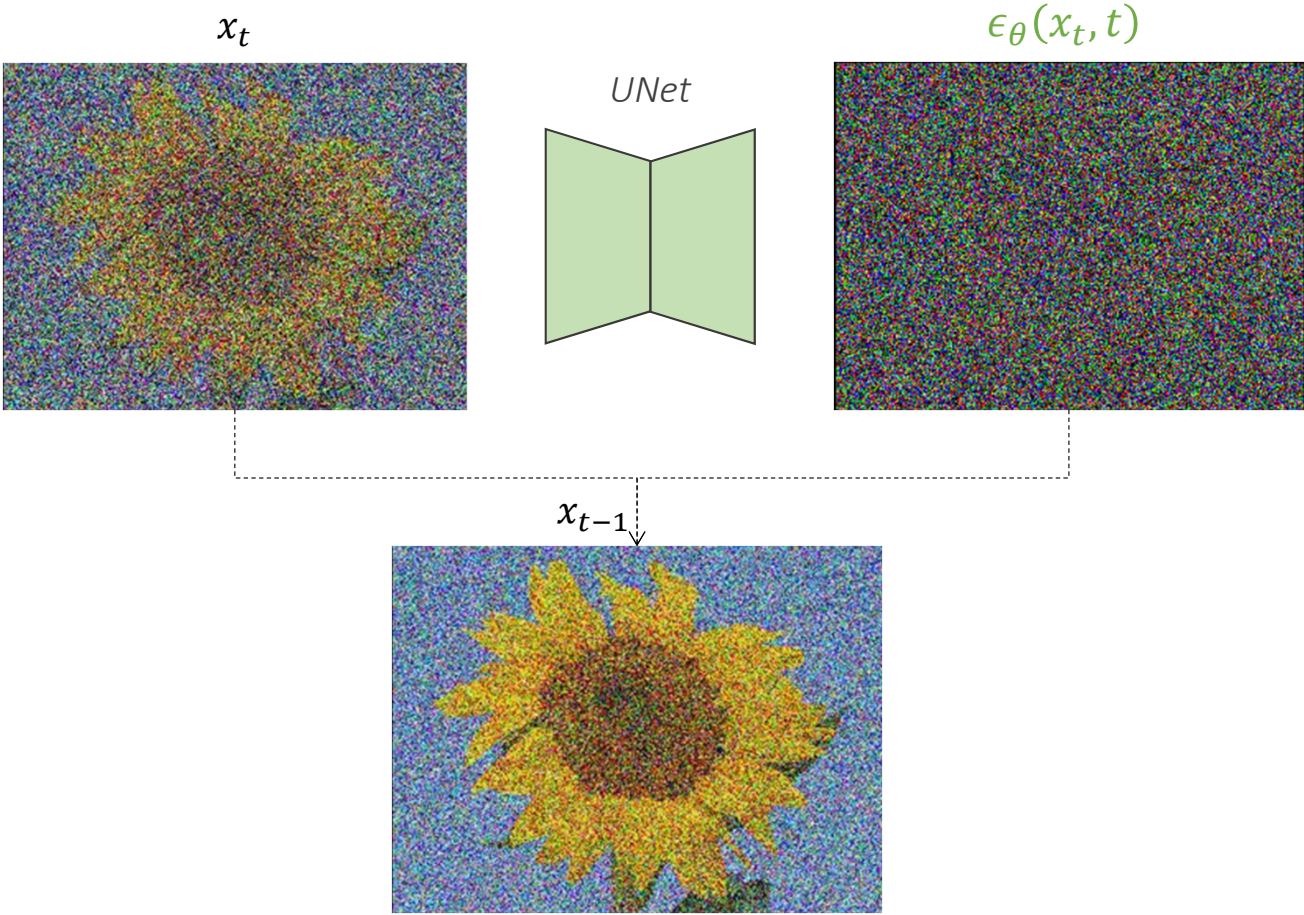
- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising



Diffusion Models

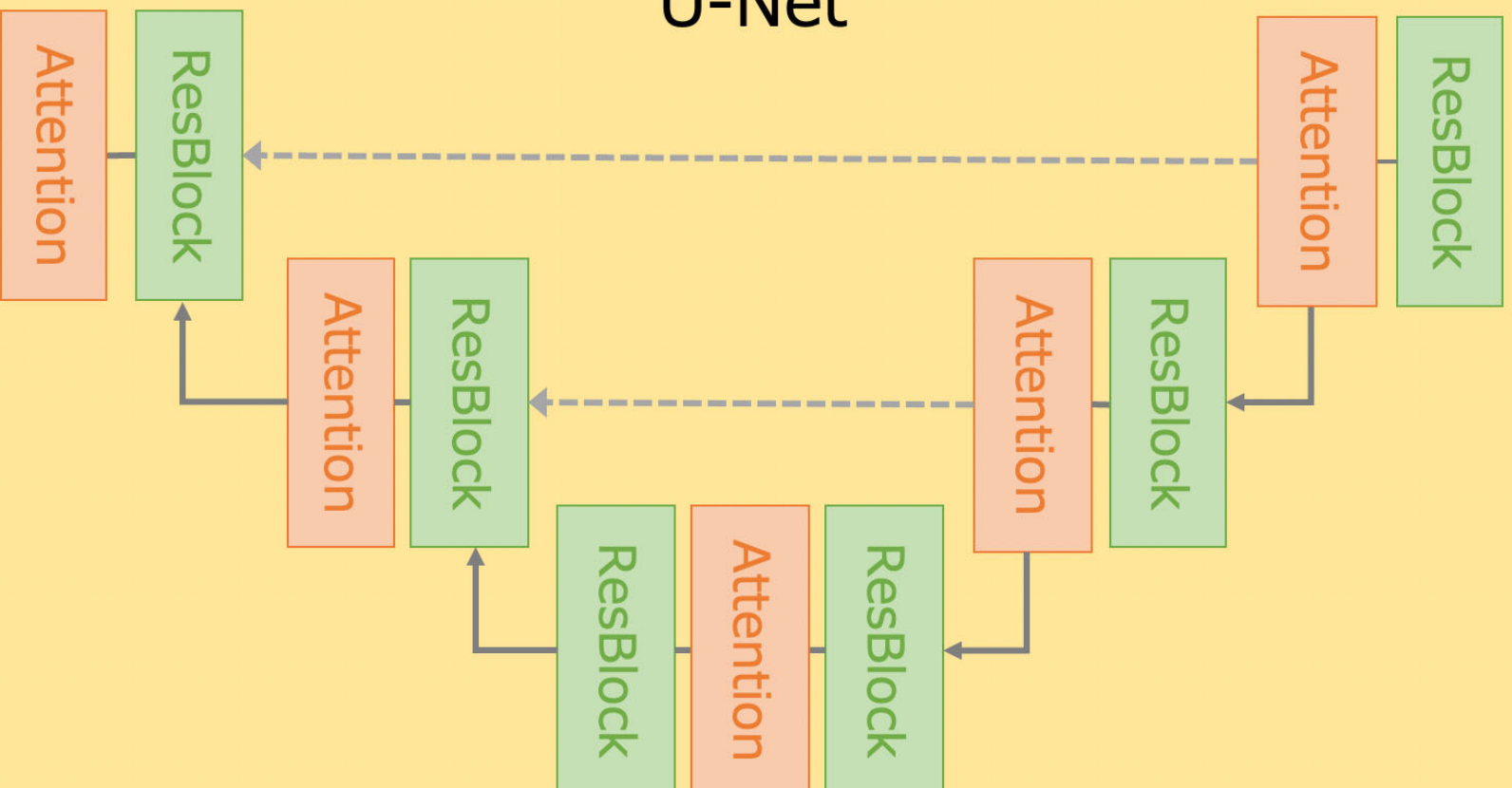


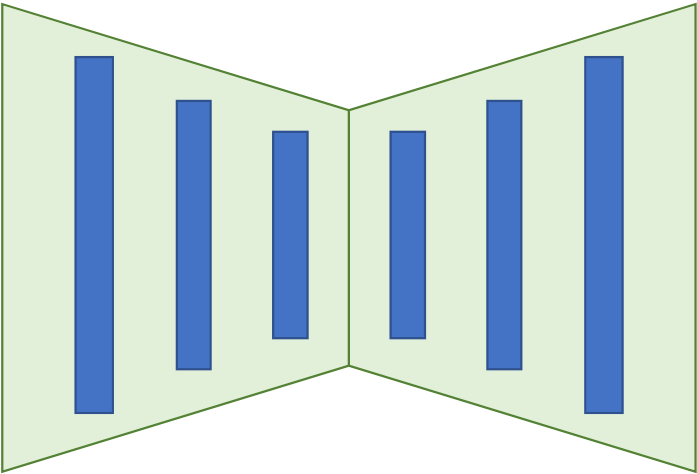
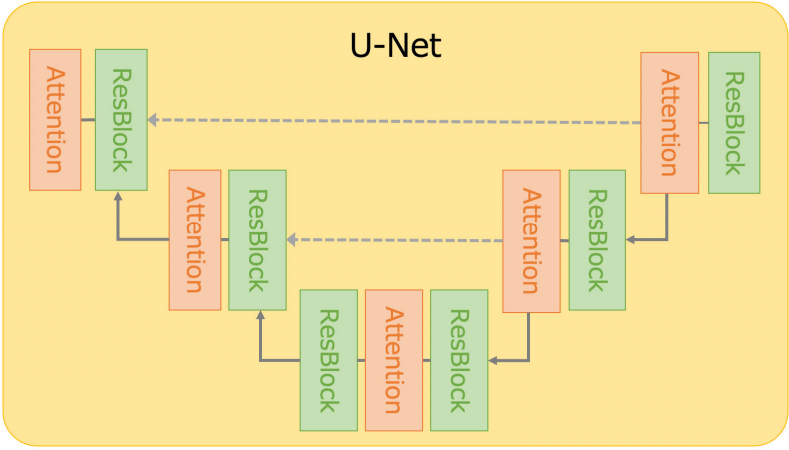
Diffusion Models



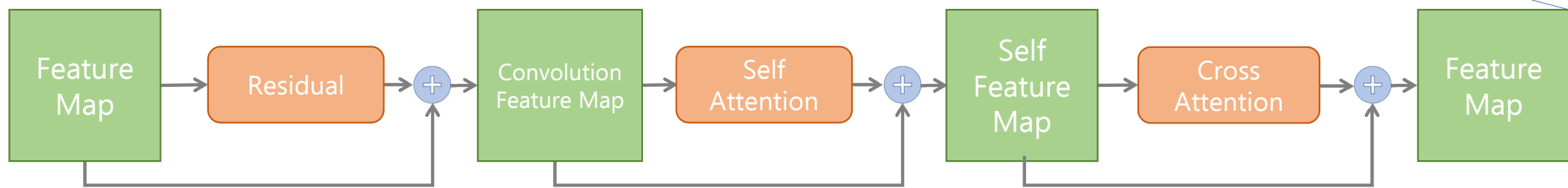
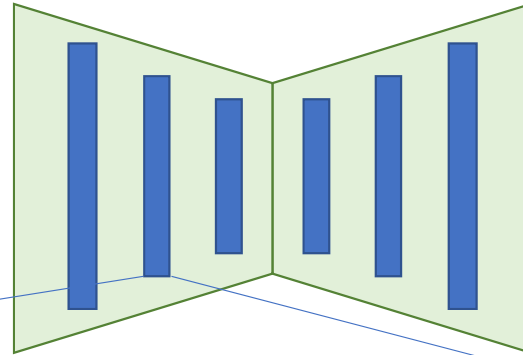
Unet

U-Net

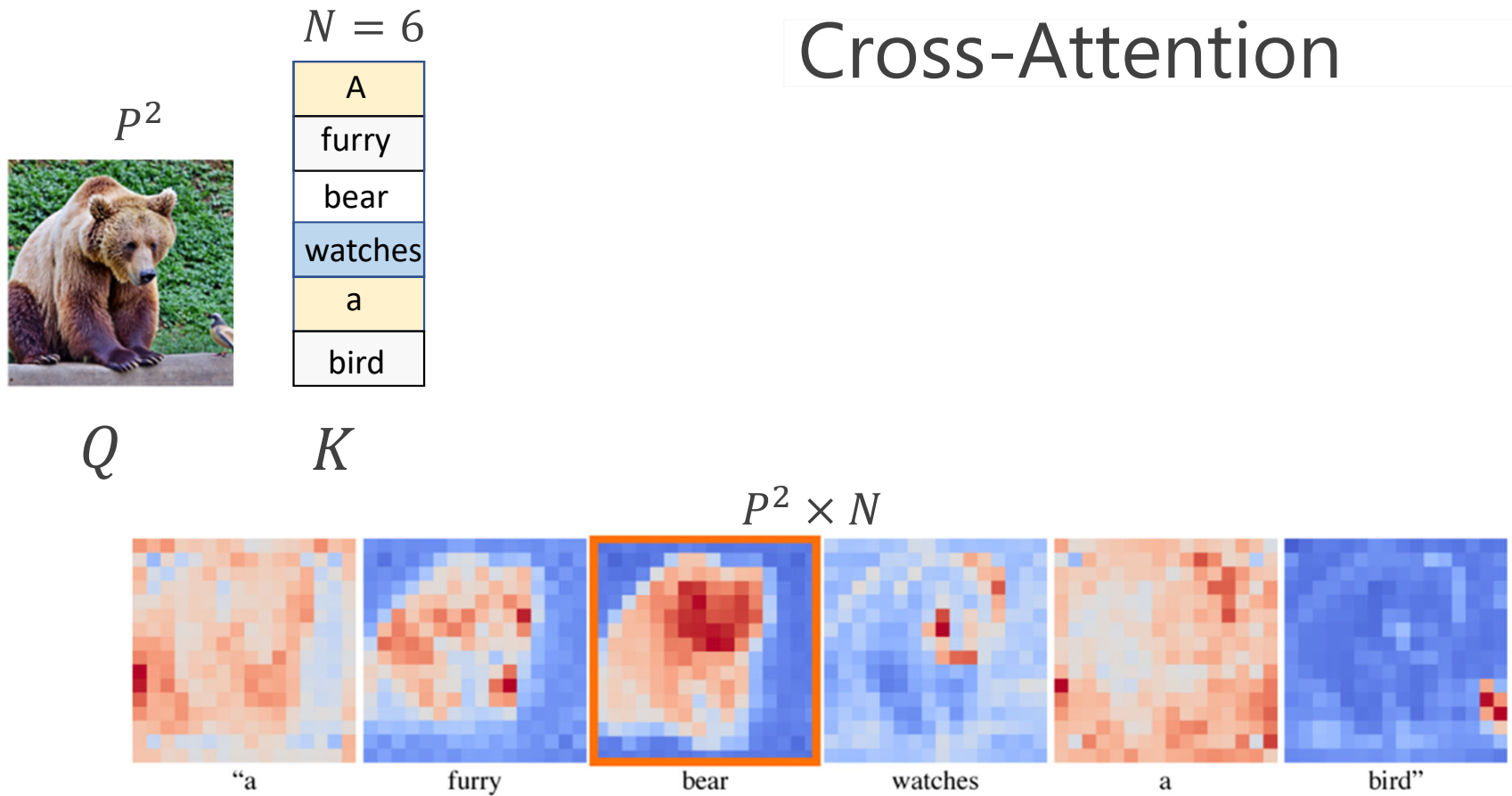




The UNet Layer



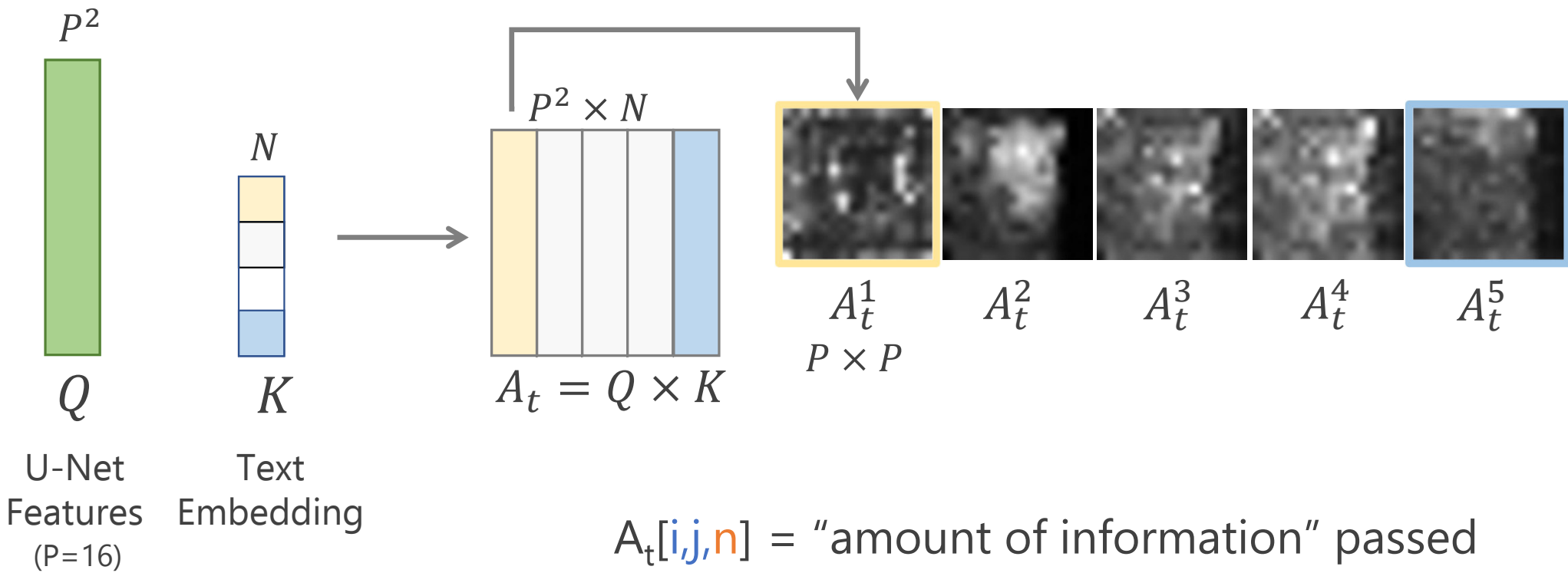
Cross-Attention



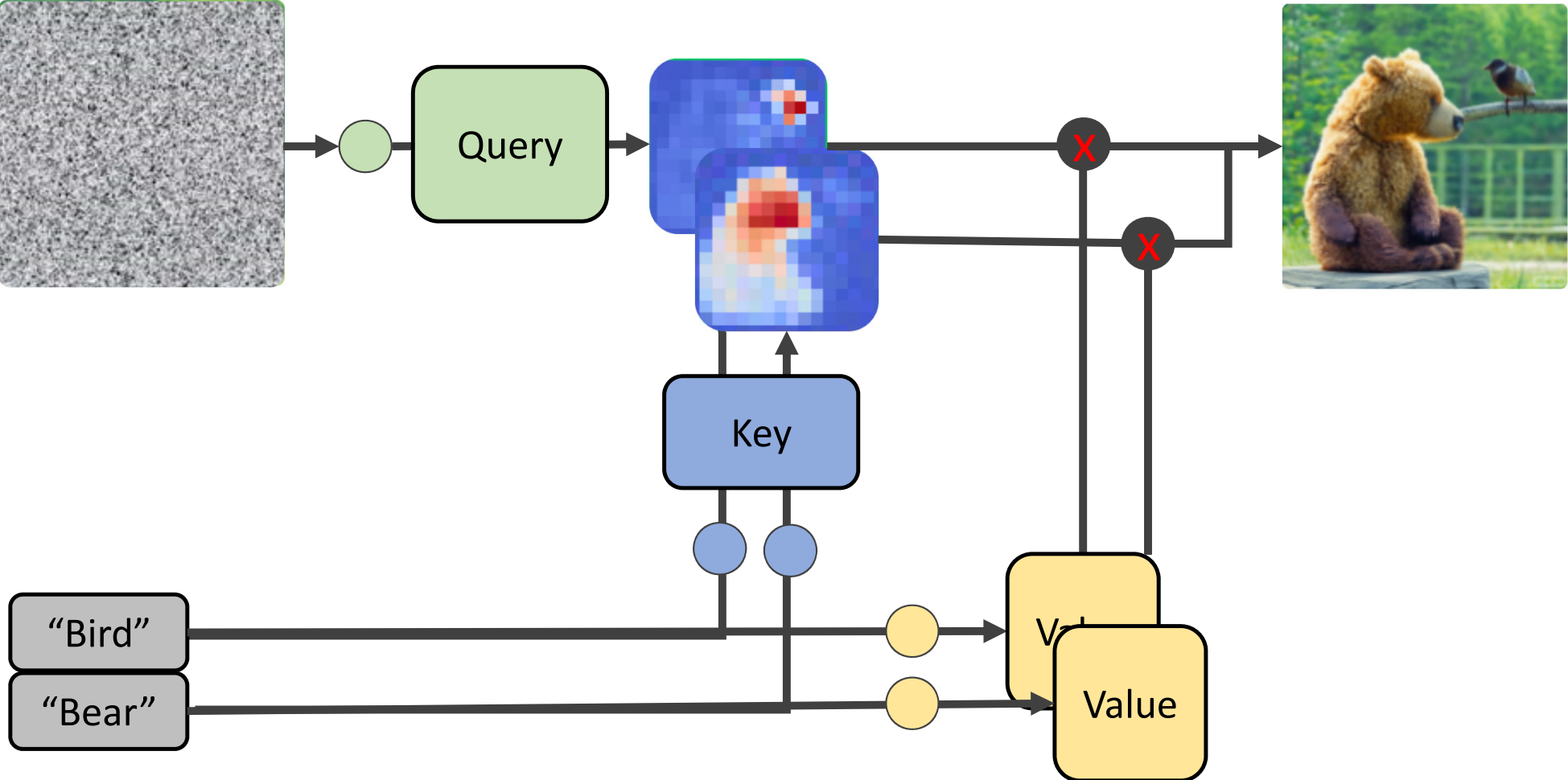
$$A_t = Q \times K$$

$A_t[i,j,n]$ = “amount of information” passed from token n to patch (i,j)

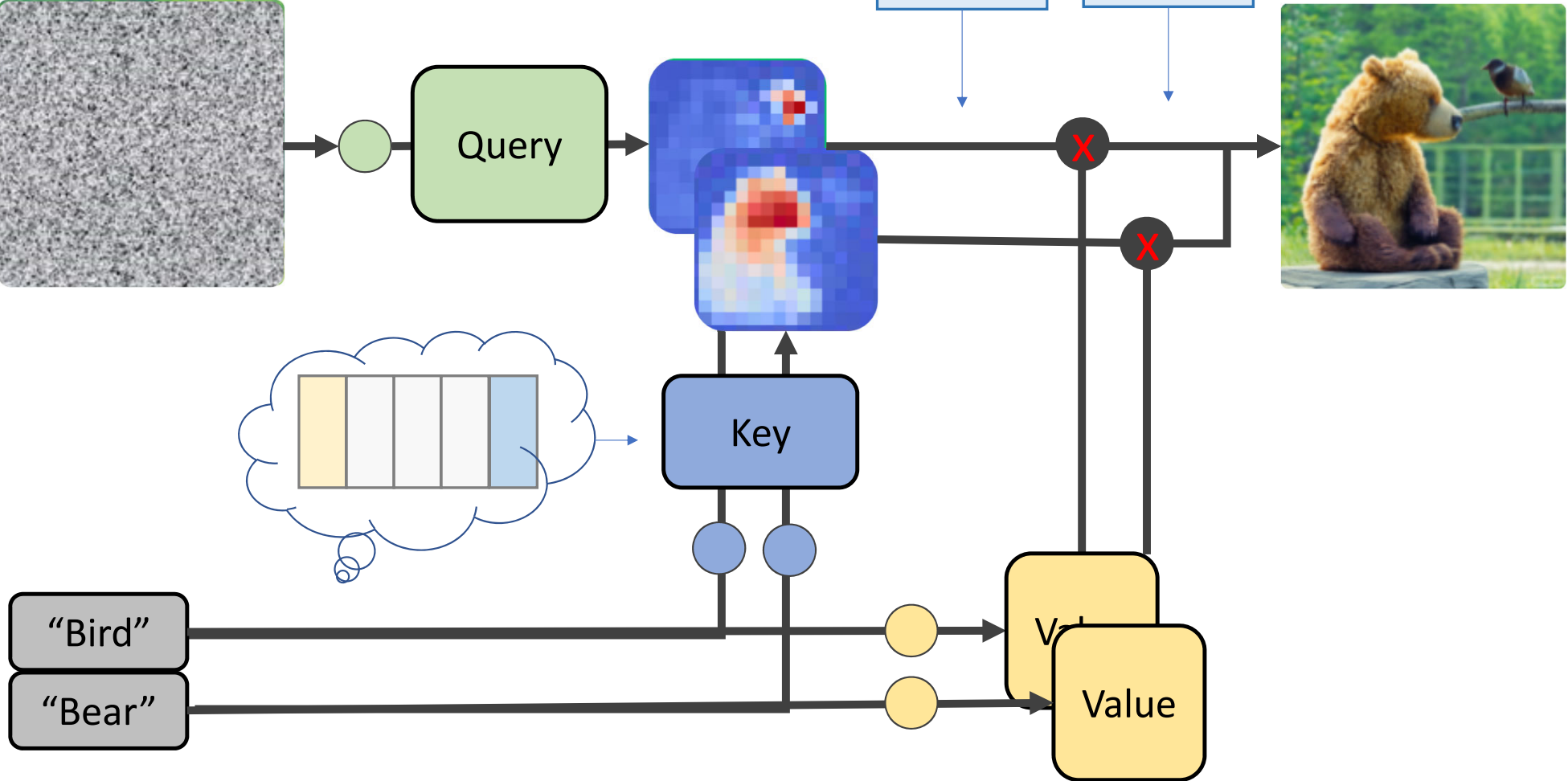
The Cross-Attention



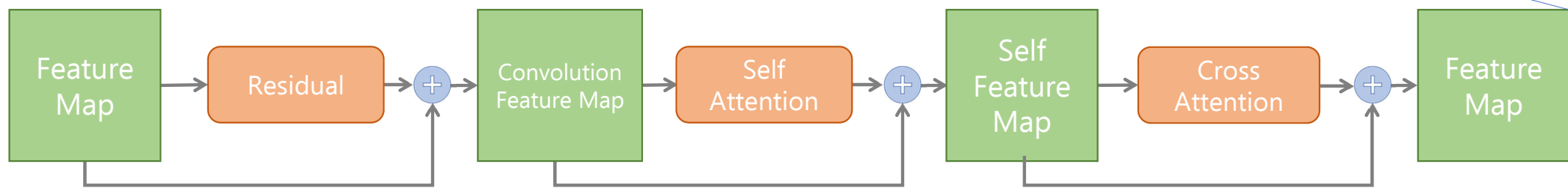
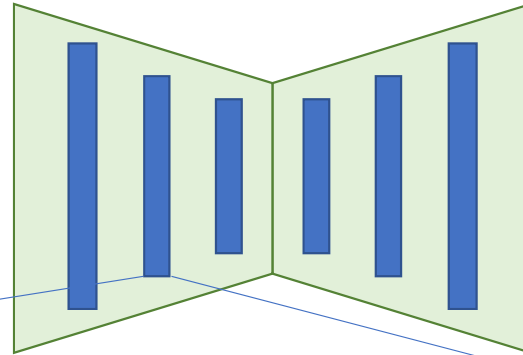
Cross-attention



Cross-attention

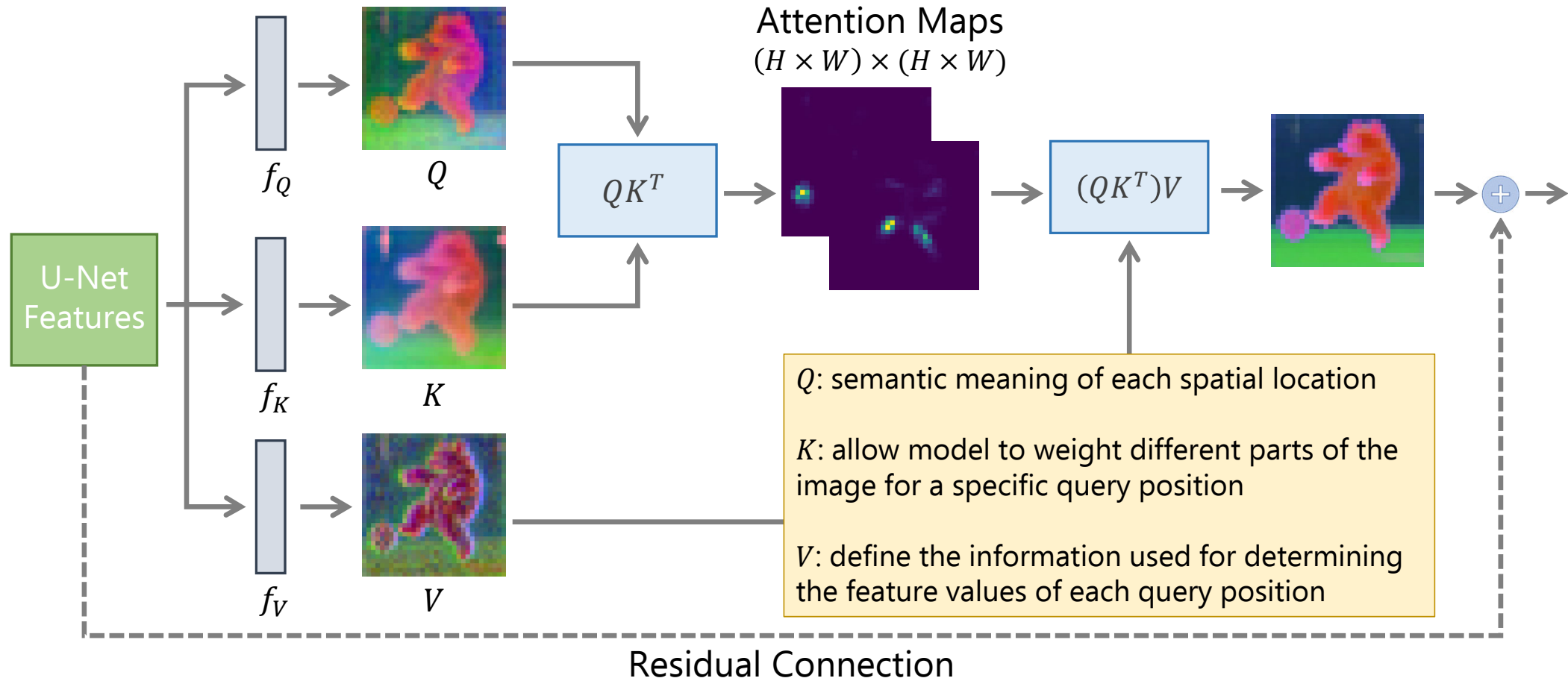


The UNet Layer

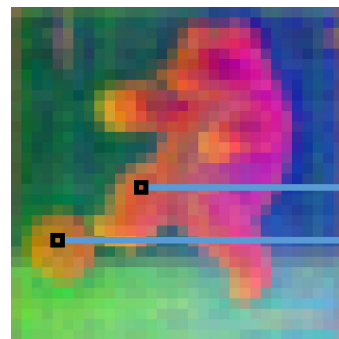


The Self-Attention

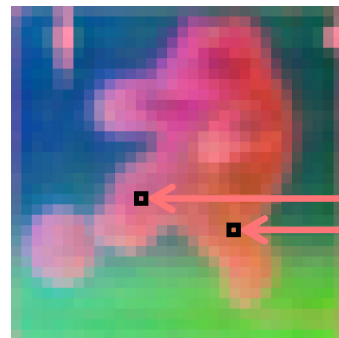
Represents where the model "looks" in the image for each spatial position in Q



The Self-Attention

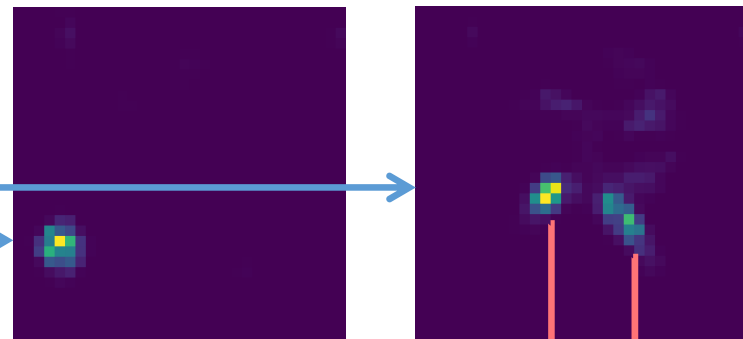


Queries
 $H \times W \times C$



Keys
 $H \times W \times C$

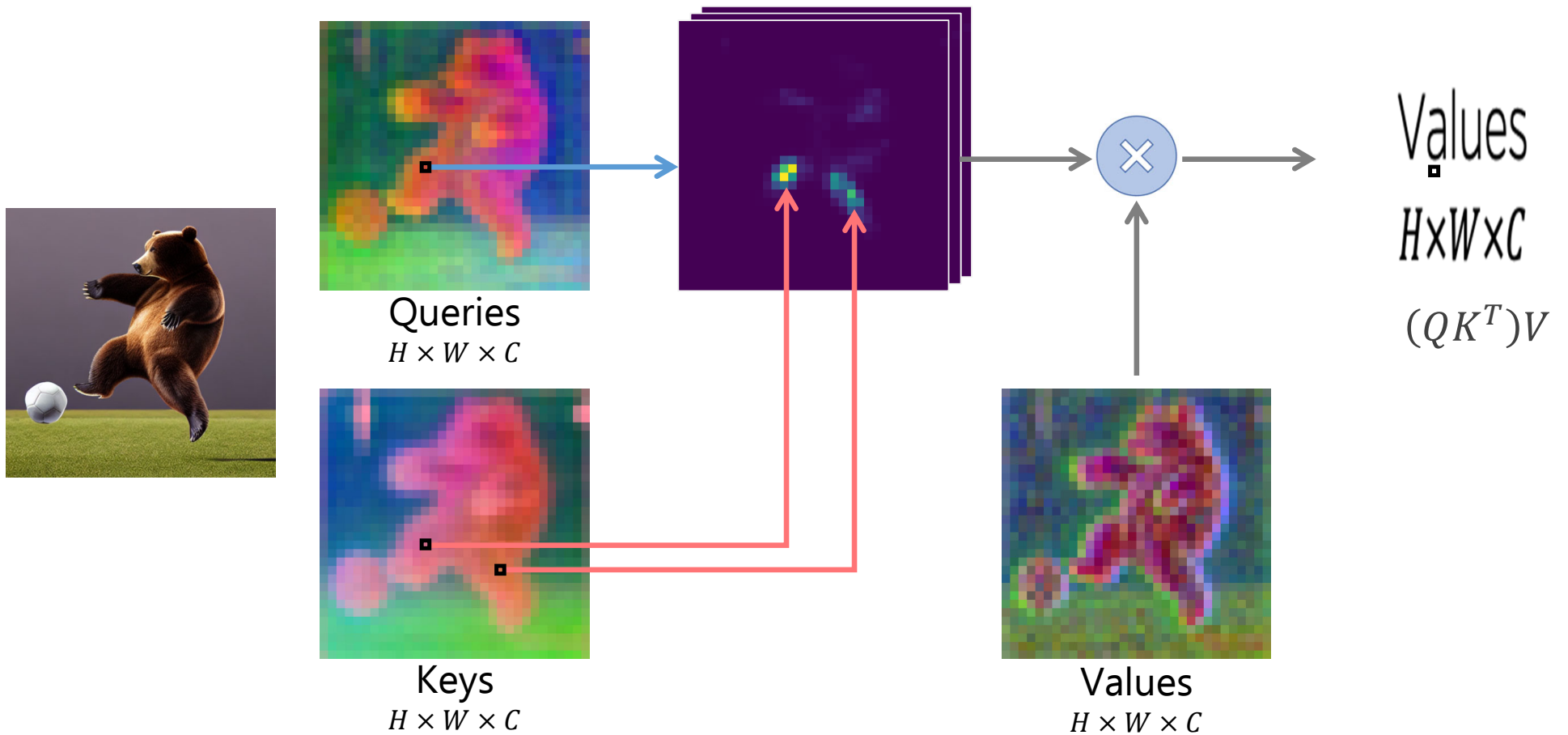
Attention Maps
 $(H \times W) \times (H \times W)$



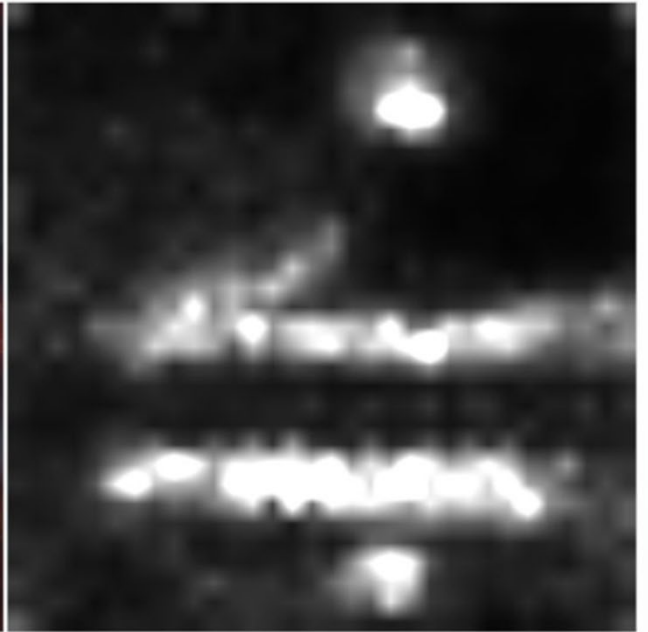
Each query defines to a
 $H \times W$ attention map!

A query on the leg of the bear "attends"
to keys located on the leg of the bear!

The Self-Attention



t = 0.6, layer: 10 / 70



t = 0.6, layer: 35 / 70



שנים (מלמעלה למטה) בורוקן 177
הורוקן 177

Self-Segmentation

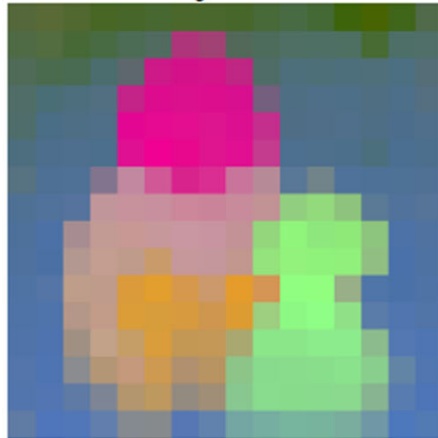
Localizing Object-level Shape Variations [Patashnik et al., ICCV 2023]

Self-Attention Maps

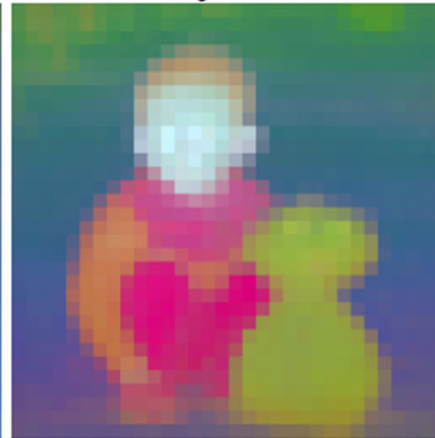
Input image



layer=4



layer=8



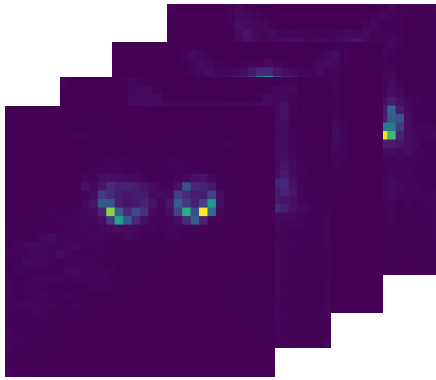
layer=11



Are these PCA on the self-attention ? On what exactly the QK maps?

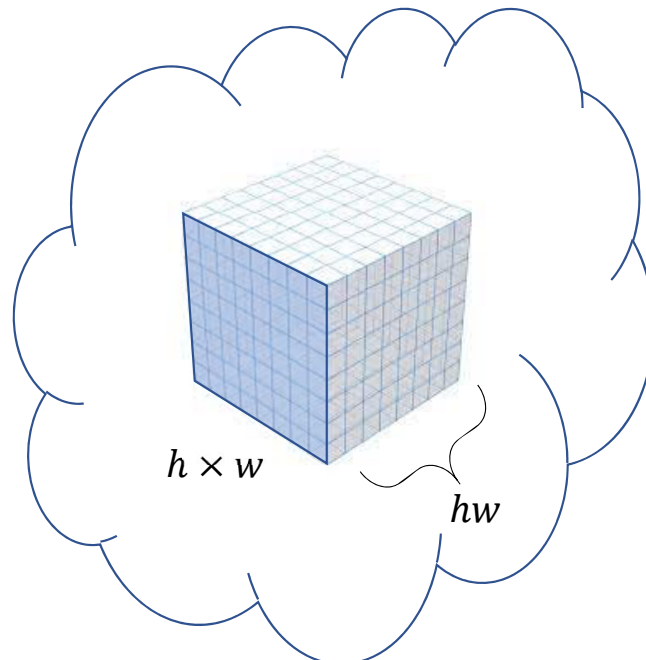
Self-Segmentation

"a **cat** is wearing sunglasses"



$hw \times (h \times w)$

There is a lot of semantics in the self attention features!!!



cluster



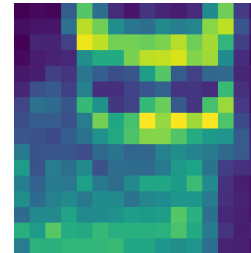
Segments labeling



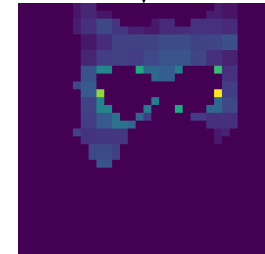
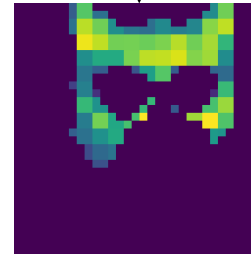
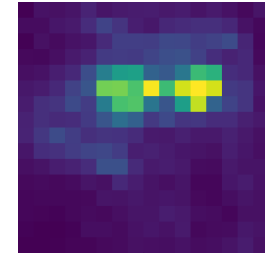
"a cat is wearing sunglasses"



cat



sunglasses



score: 0.65

score: 0.19

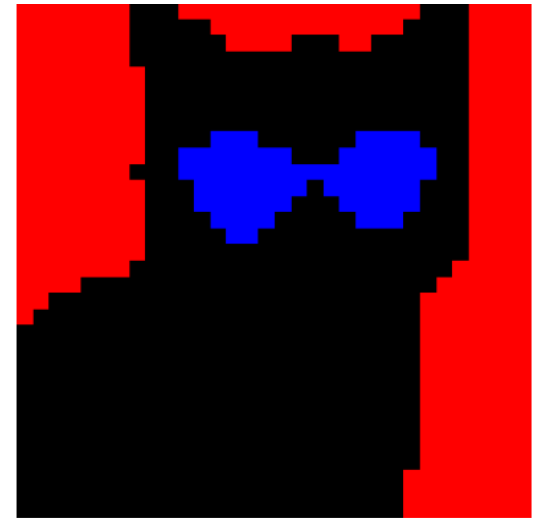
Segments labeling



"a cat is wearing sunglasses"



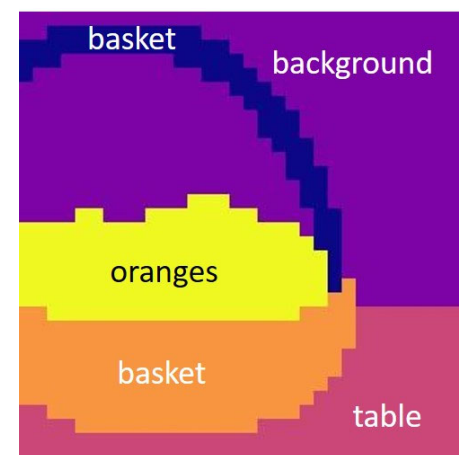
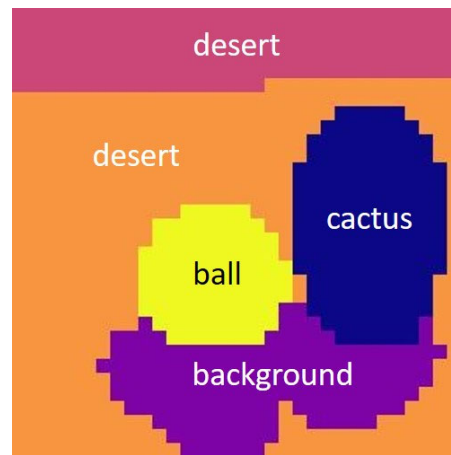
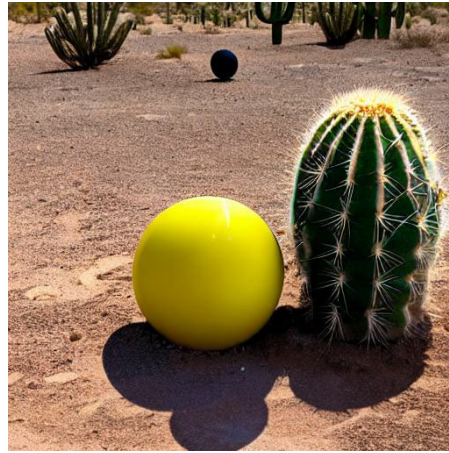
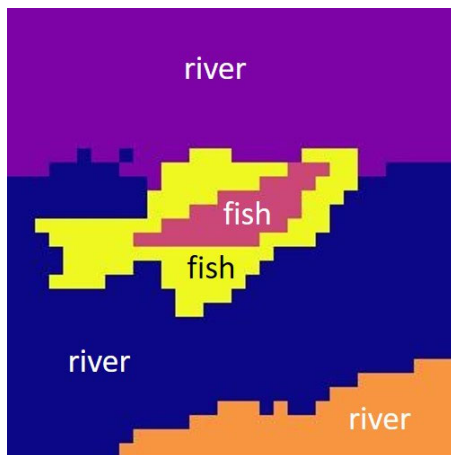
■ 1	■ 1	■ 4
■ background	■ background	■ 1



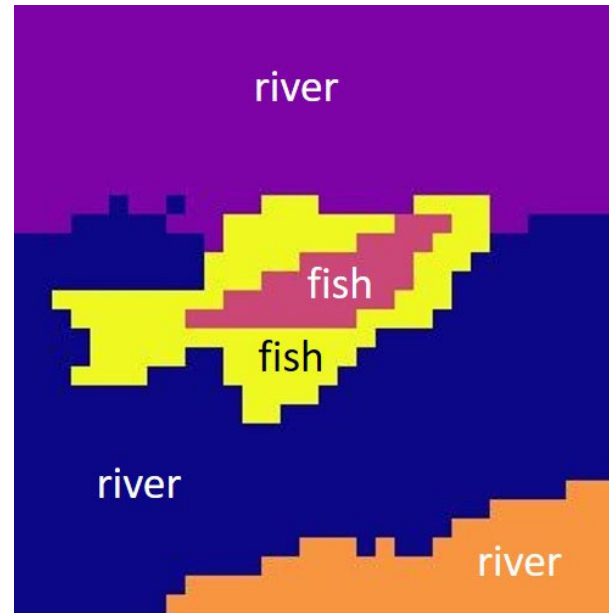
■ 1	■ background	■ 4
-----	--------------	-----

1-cat, 4-sunglasses

Self-Segmentation Results



Self-Segmentation Results



Cross-Image Attention

Cross-Image Attention for Zero-Shot Appearance Transfer

[Yuval Alaluf](#)^{*}, Daniel Garibi^{*}, Or Patashnik, [Hadar Averbuch Elor](#),
[Daniel Cohen-Or](#)



Motivation



Structure

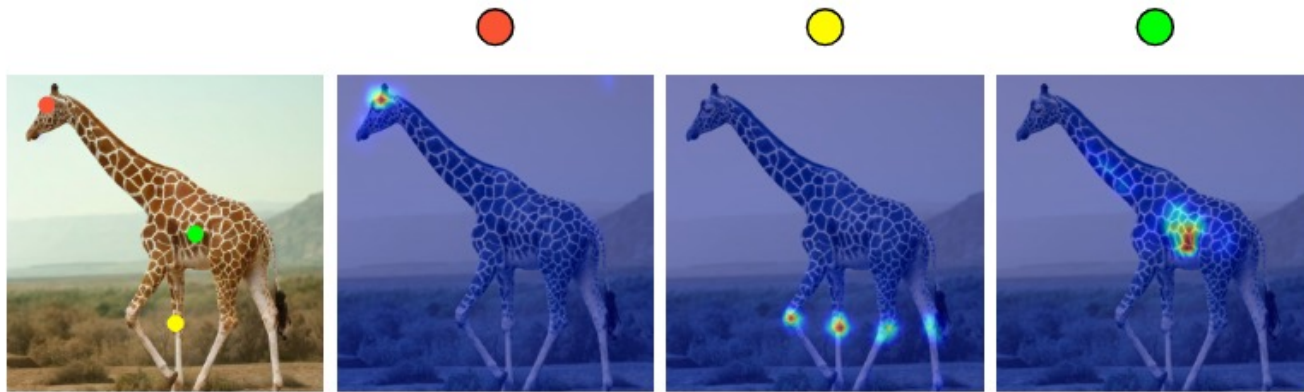


Appearance

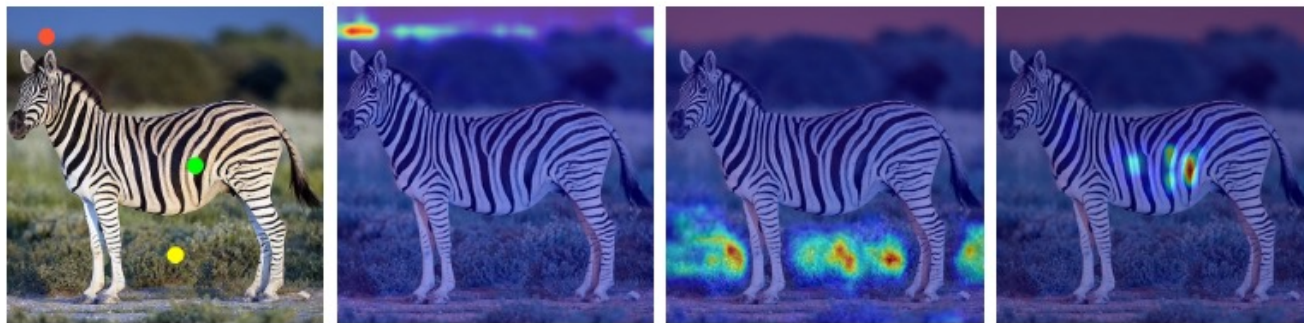


Output

The Roles of the Queries, Keys, and Values



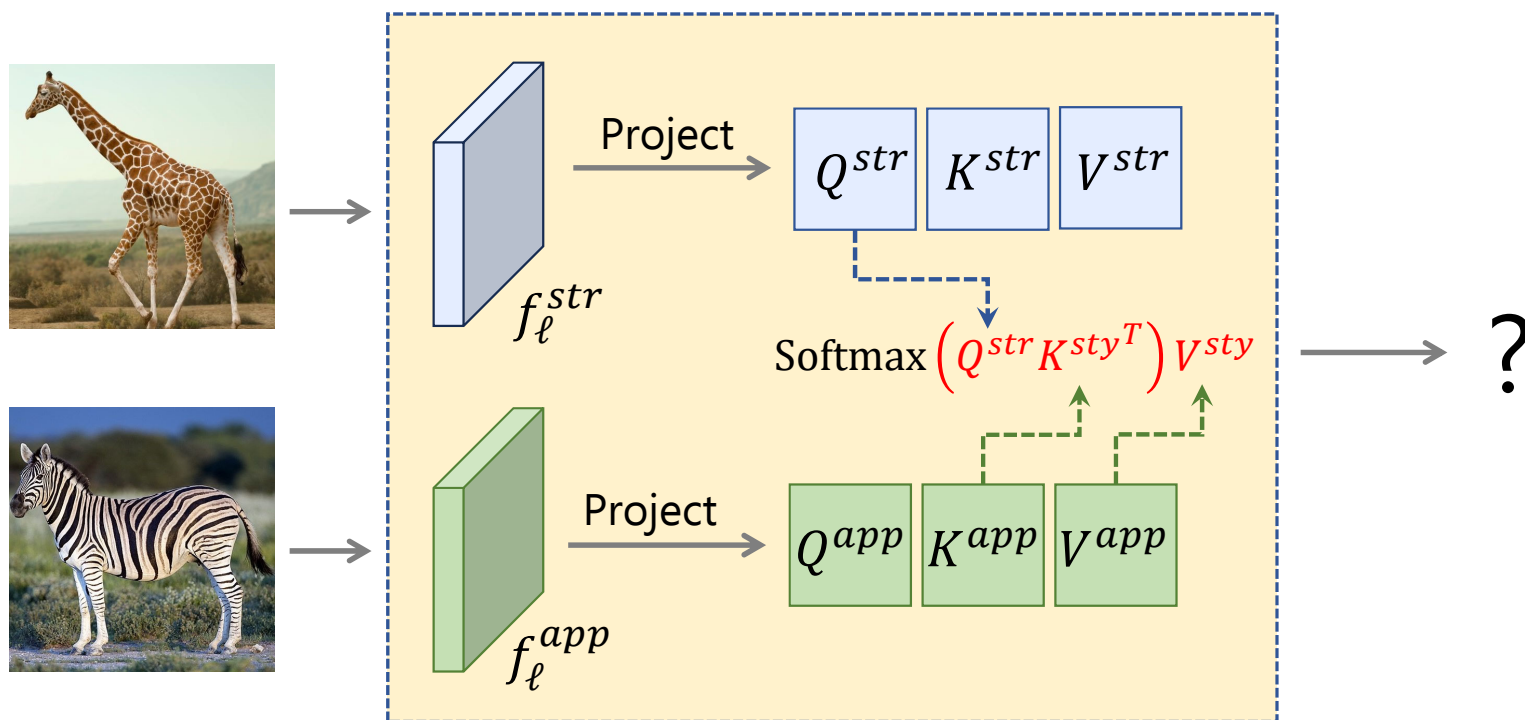
$$Q_{struct} \cdot K_{struct}^T$$



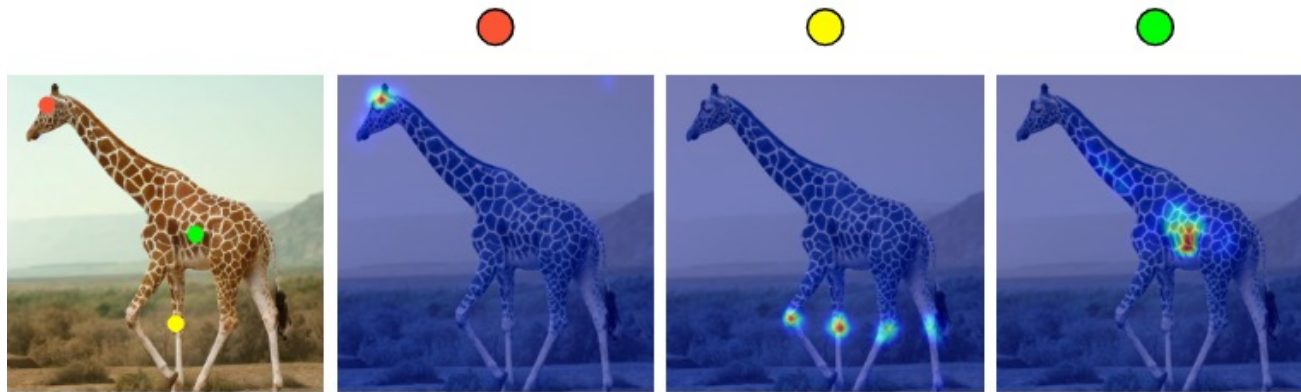
$$Q_{app} \cdot K_{app}^T$$

Self-attention maps, which focus on semantically similar regions in the image.

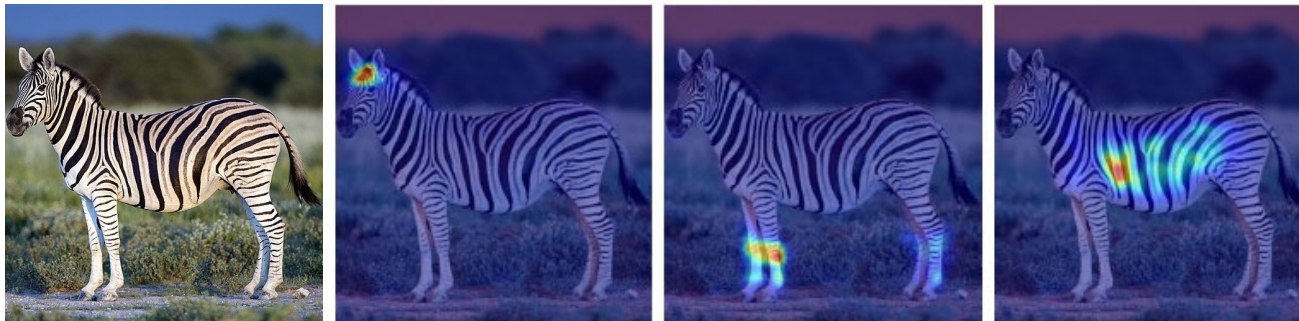
What If we Swapped the Queries, Keys, and Values Between Different Images?



The Roles of the Queries, Keys, and Values



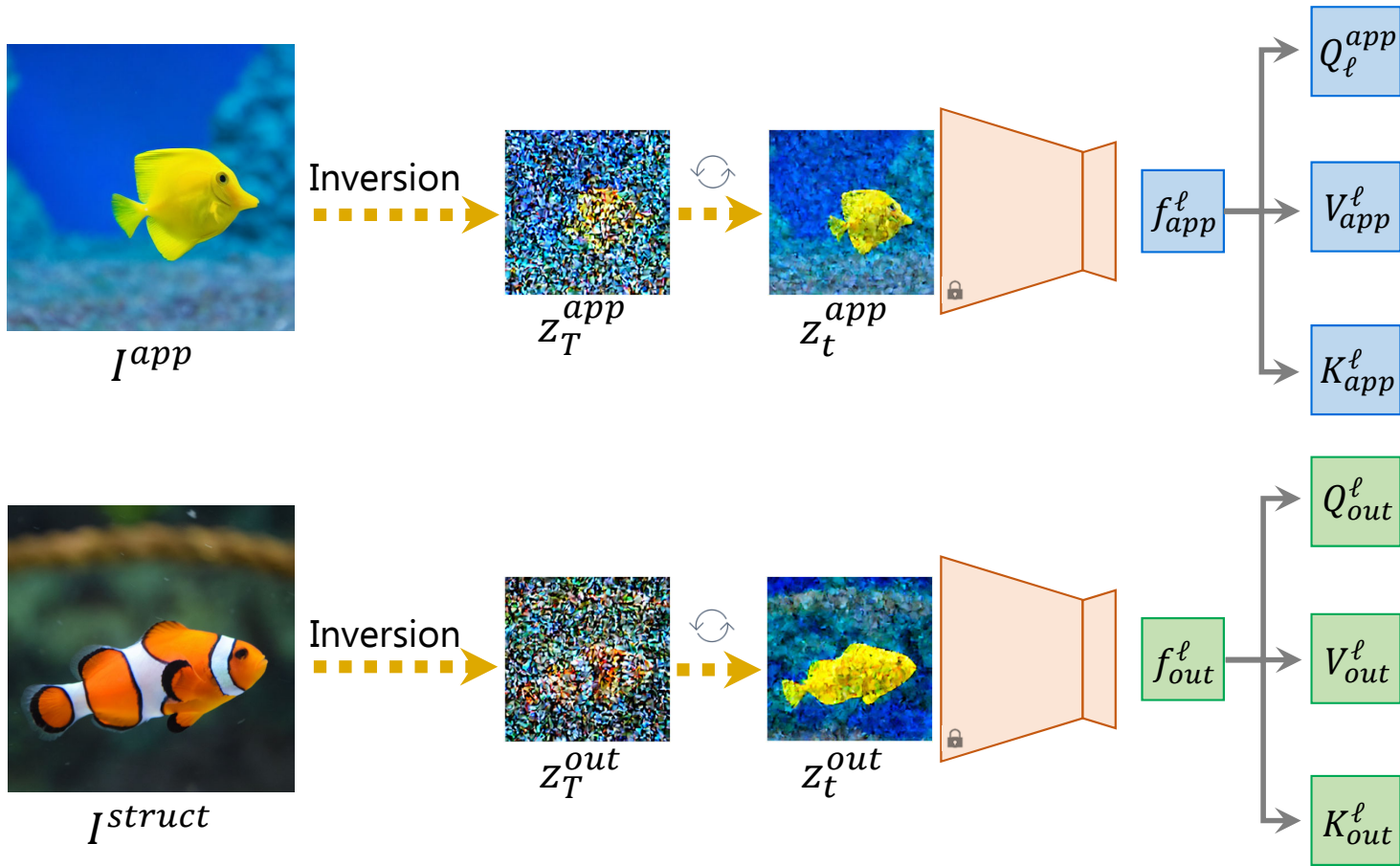
$$Q_{struct} \cdot K_{struct}^T$$



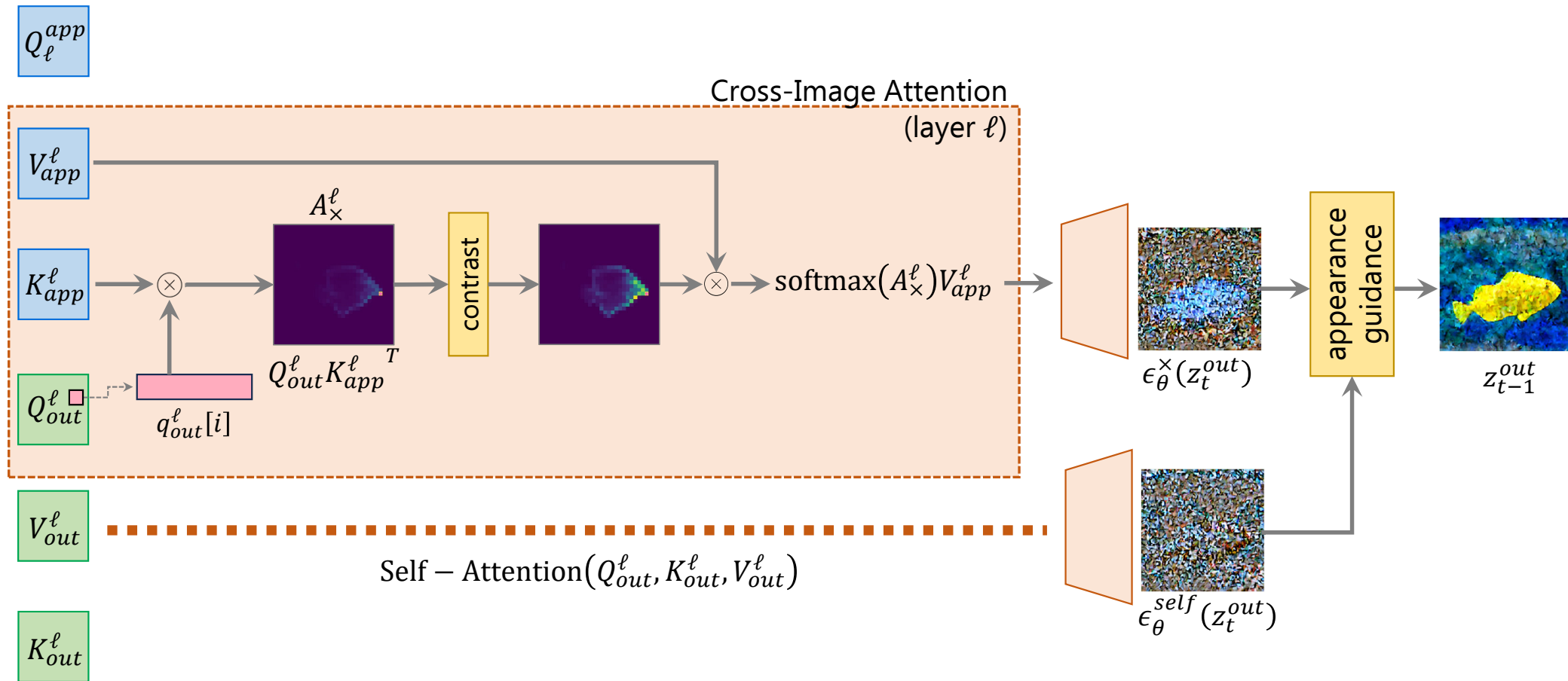
$$Q_{struct} \cdot K_{app}^T$$

Taking the **queries** from the structure image and the **keys** from the appearance image gives semantic correspondences between objects!

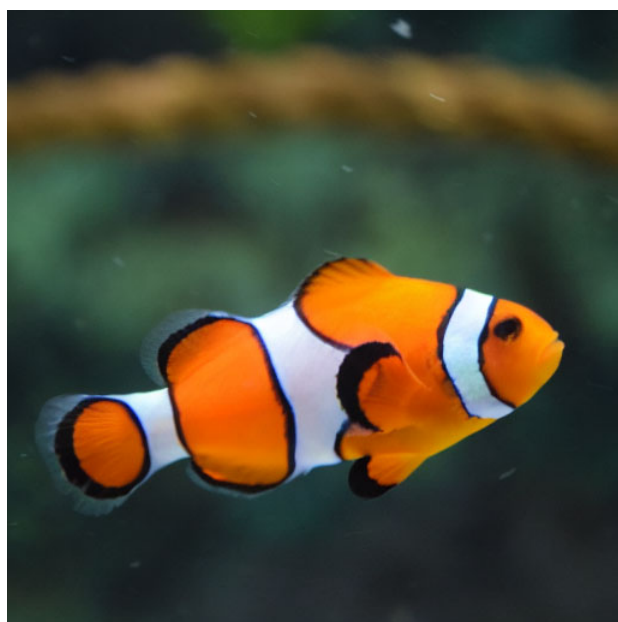
The Cross-Image Attention



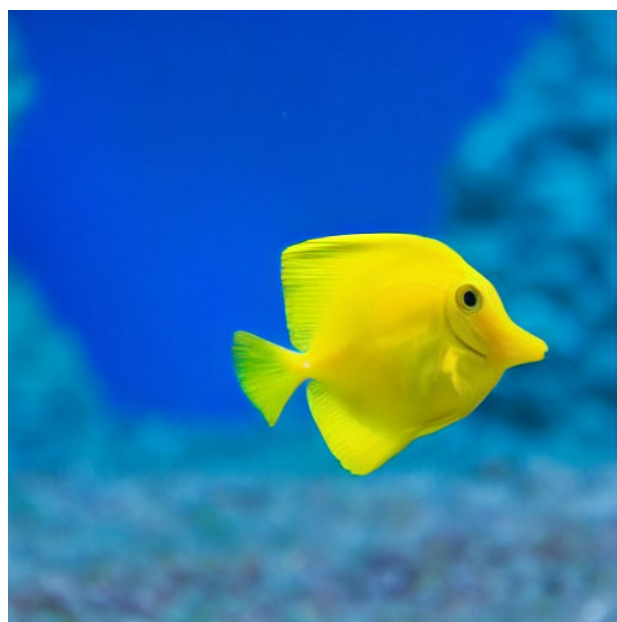
The Cross-Image Attention



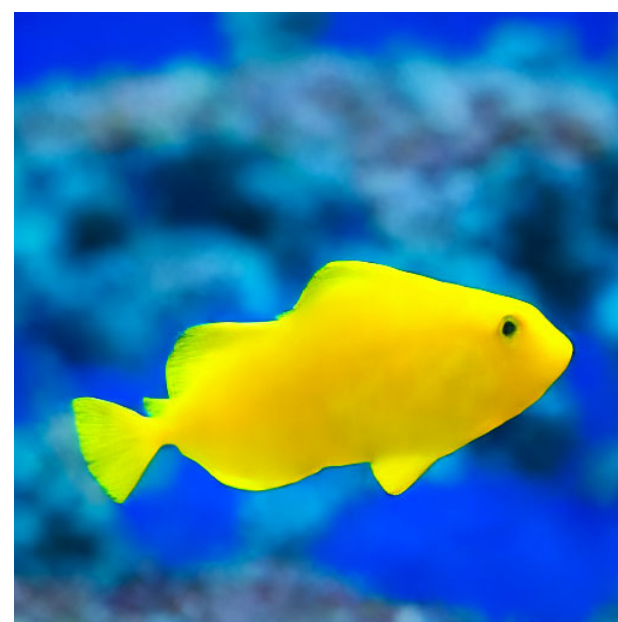
The Cross-Image Attention



Structure



Appearance



Output

Appearance Transfer Results



Structure



Appearance



Output

Appearance Transfer Results



Structure



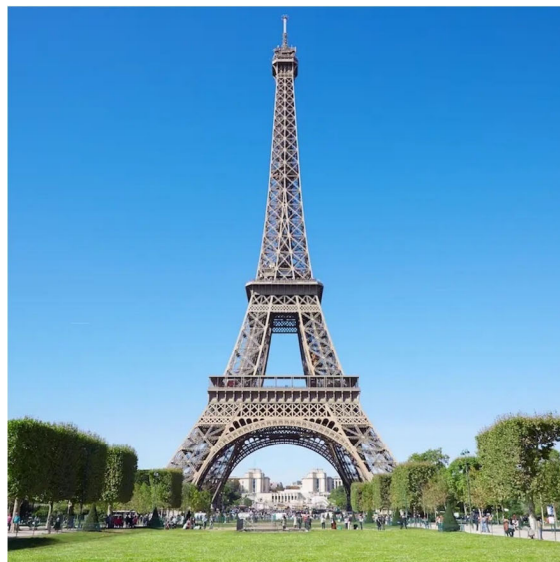
Appearance



Output

Appearance Transfer Results

Eiffel Tower



Structure

Sagrada Família



Appearance



Output

Appearance Transfer Results



Structure



Appearance



Output

Appearance Transfer Results





Structure



Appearance



Output

StyleAlign

[\(link\)](#)

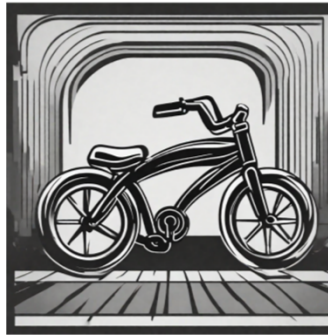
“Toy train...”



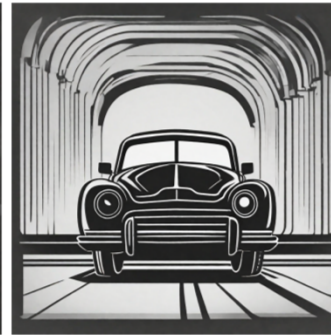
“Toy airplane...”



“Toy bicycle...”



“Toy car...”



“Toy boat...”



“...BW logo, high contrast.”

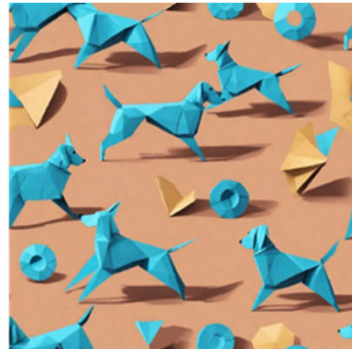


“...colorful, macro photo.”

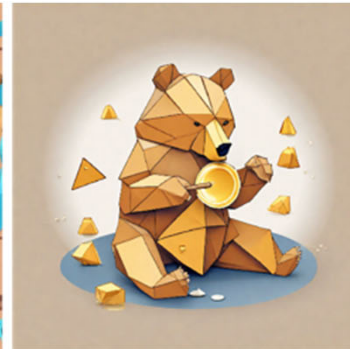
Text-to-Image Generation



"A cat playing with a ball of wool..."



"A dog catching a frisbee..."



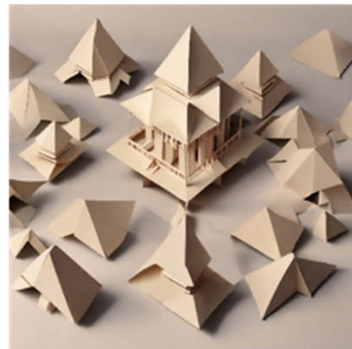
"A bear eating honey..."



"A whale playing with a ball..."



"A woman working in the office..."



"A temple..."



"A person riding a bike..."



"A cactus..."

"... in minimal origami style."

Text-to-Image Generation with StyleAligned



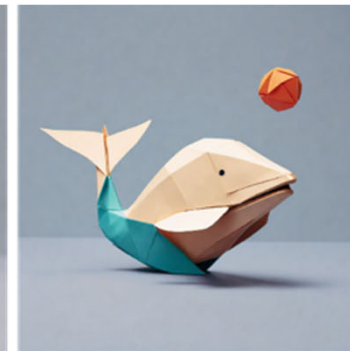
"A cat playing with a ball of wool..."



"A dog catching a frisbee..."



"A bear eating honey..."



"A whale playing with a ball..."



"A woman working in the office..."



"A temple..."



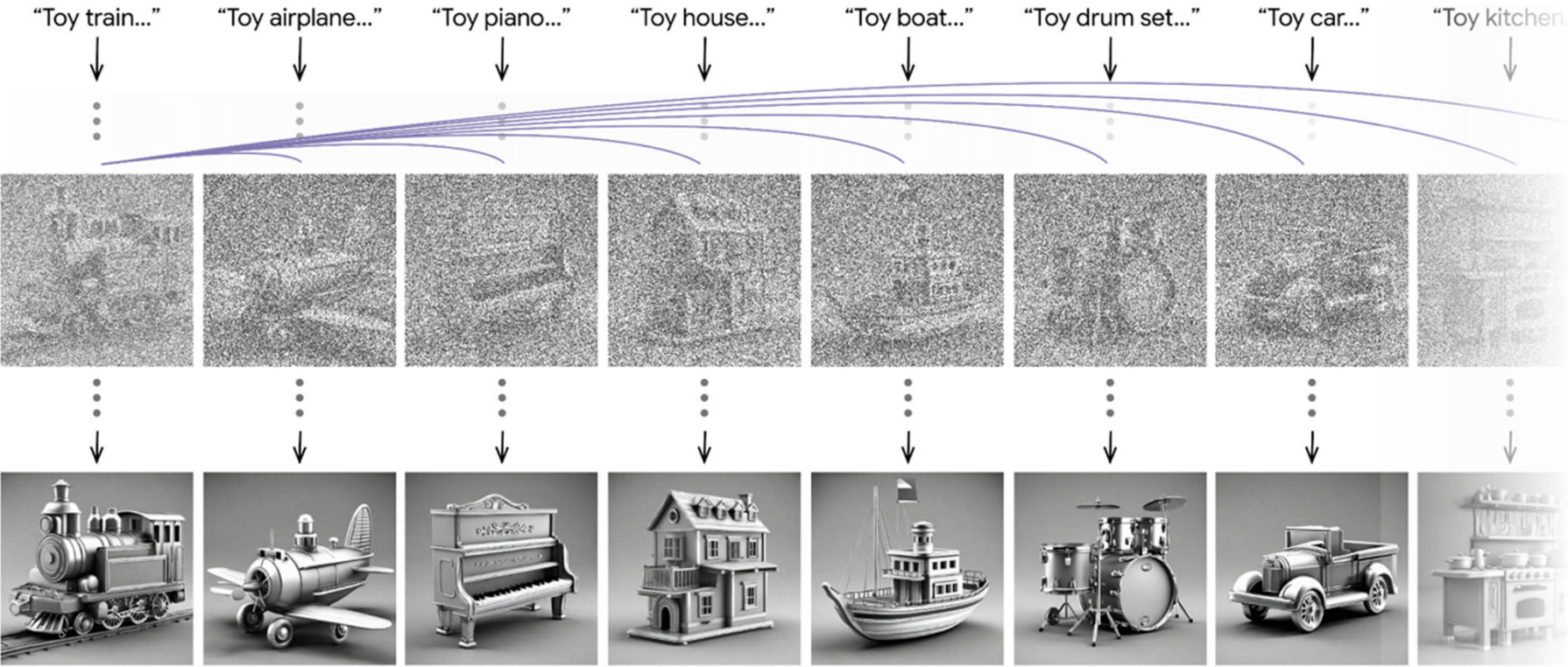
"A person riding a bike..."



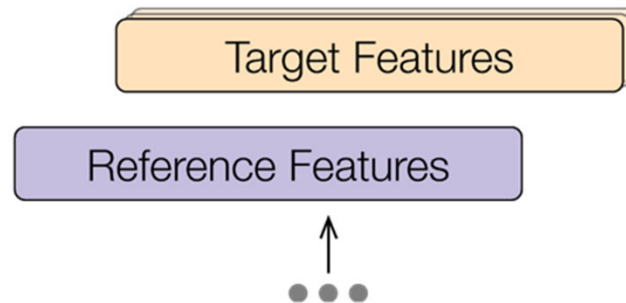
"A cactus..."

"... in minimal origami style."

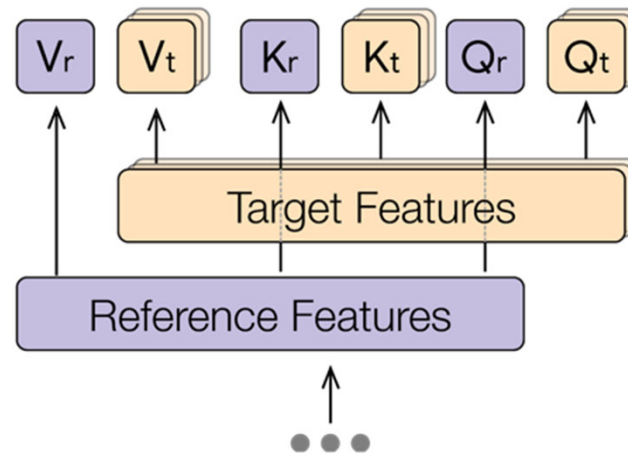
Shared attention during the diffusion process



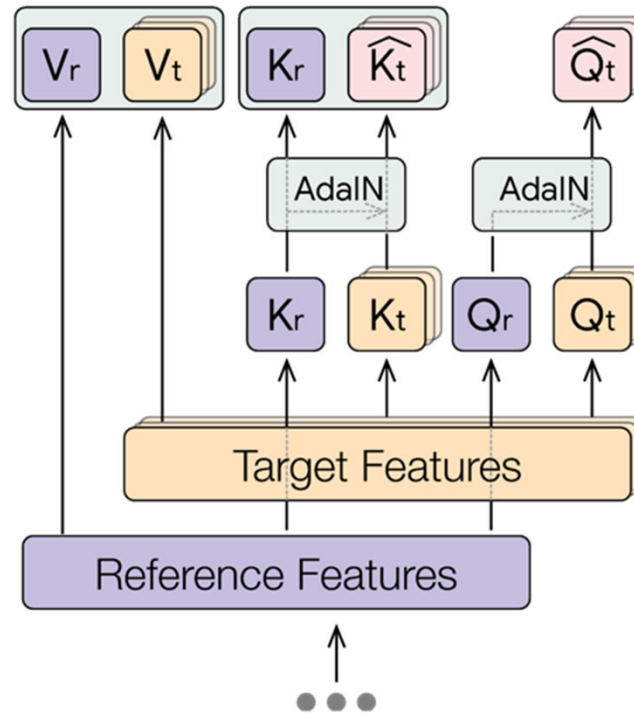
Shared Attention Layer



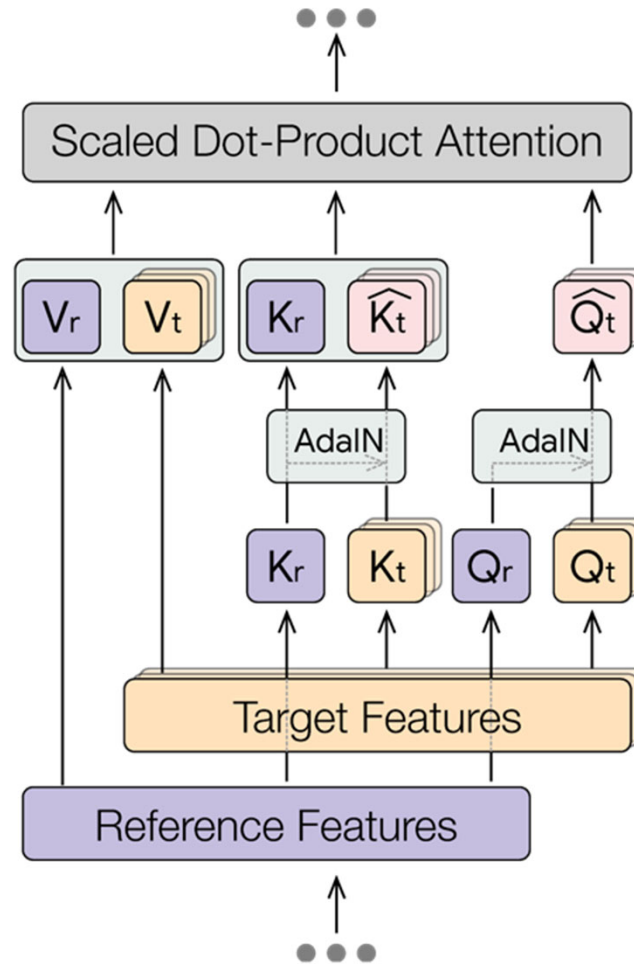
Shared Attention Layer



Shared Attention Layer



Shared Attention Layer



Style Aligned generation of Synthetic Images

“Firewoman...”



“Gardner...”



“Scientist ...”



“Police woman...”



“Saxophone player...”



“Painter...”



“Astronaut...”



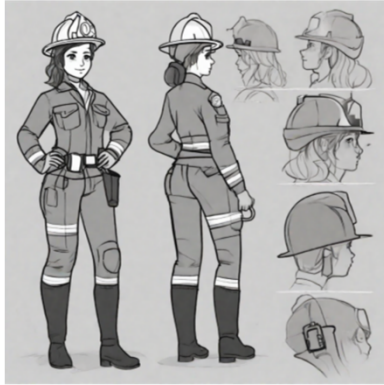
“Taxi Driver...”



“...made of claymation, stop motion animation.”

Style Aligned generation of Synthetic Images

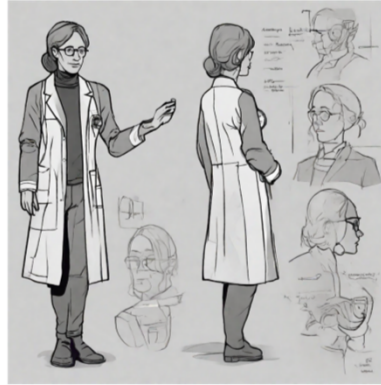
“Firewoman...”



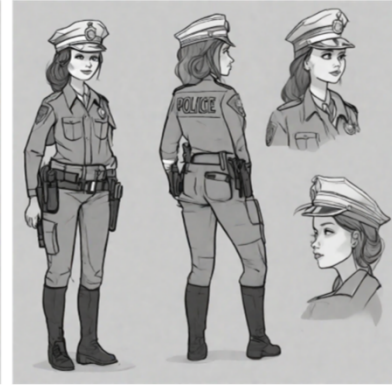
“Gardner...”



“Scientist ...”



“Police woman...”



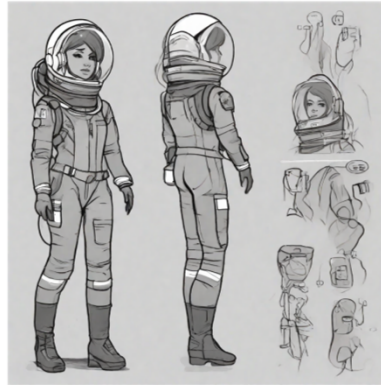
“Saxophone player...”



“Painter...”



“Astronaut...”



“Taxi Driver...”



“...sketch, character sheet.”

Style Aligned generation of Synthetic Images

“Firewoman...”



“Gardner...”



“Scientist ...”



“Police woman...”



“Saxophone player...”



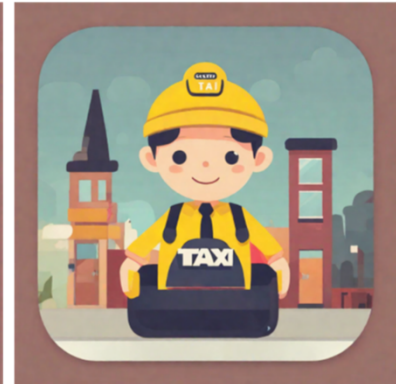
“Painter...”



“Astronaut...”



“Taxi Driver...”



“ ...in minimal flat design illustration.”

Style Aligned generation from an Input Image

Reference image



Space rocket



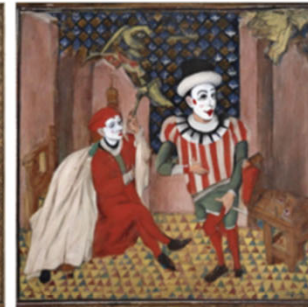
Boy riding a bicycle



Matterhorn mountain



Mime artist



Seattle needle



Style Aligned generation from an Input Image

Reference image



Space rocket



Boy riding a bicycle



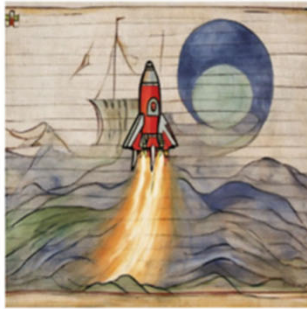
Matterhorn mountain



Mime artist

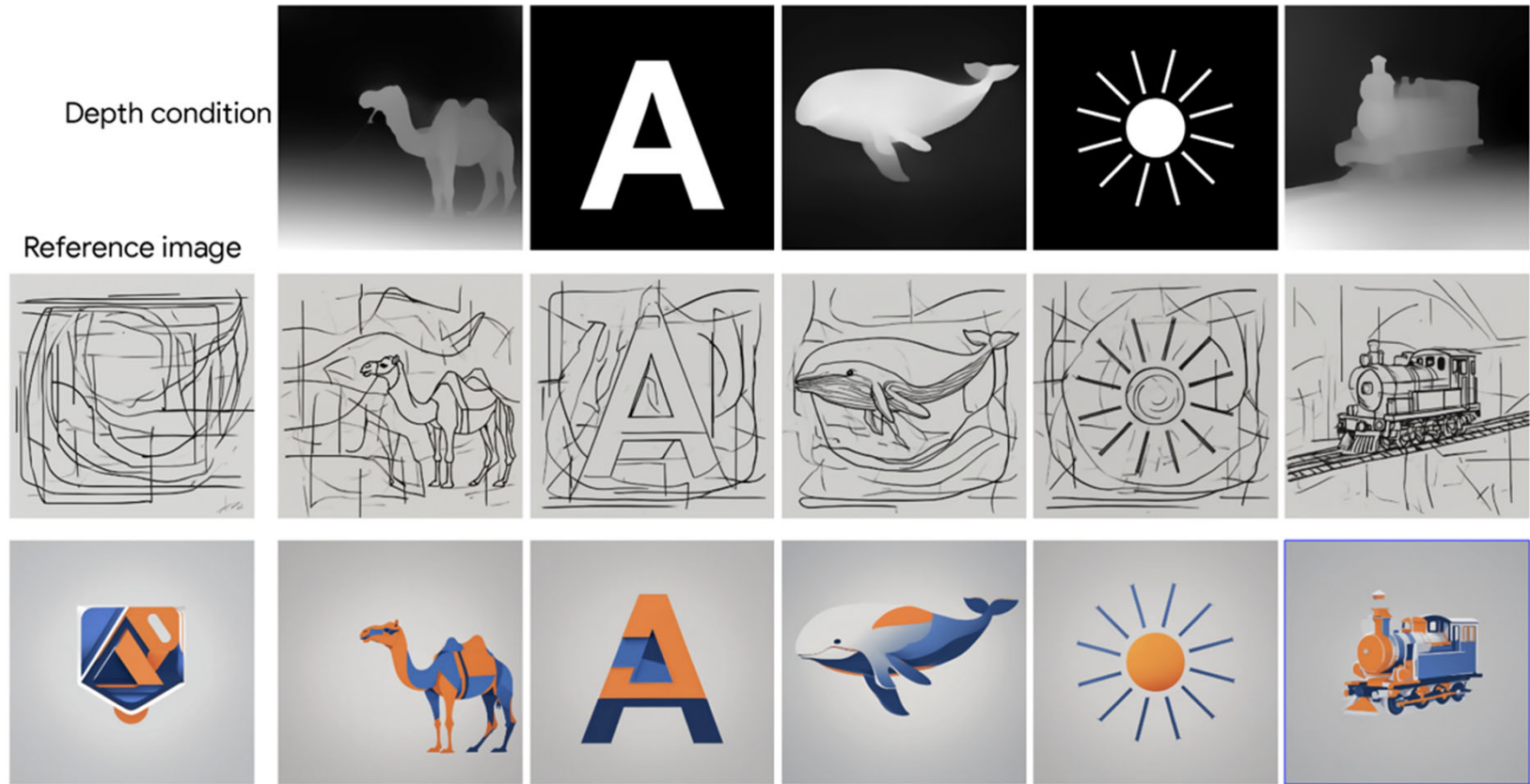


Seattle needle



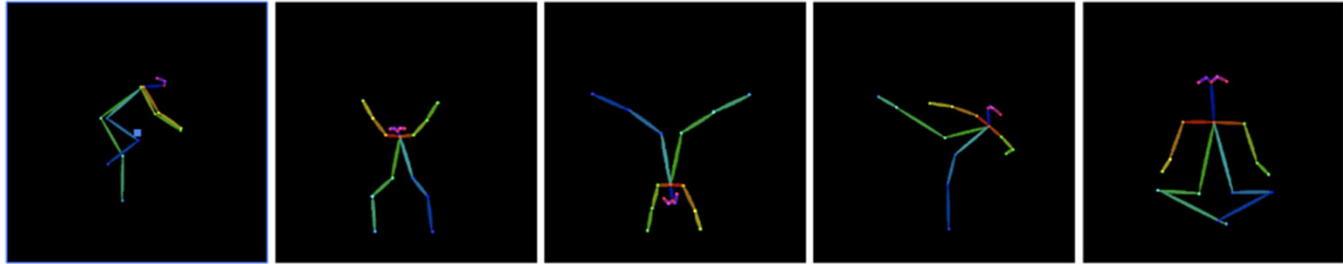
StyleAligned with other methods

ControlNet + StyleAligned



ControlNet + StyleAligned

Pose condition



Reference image

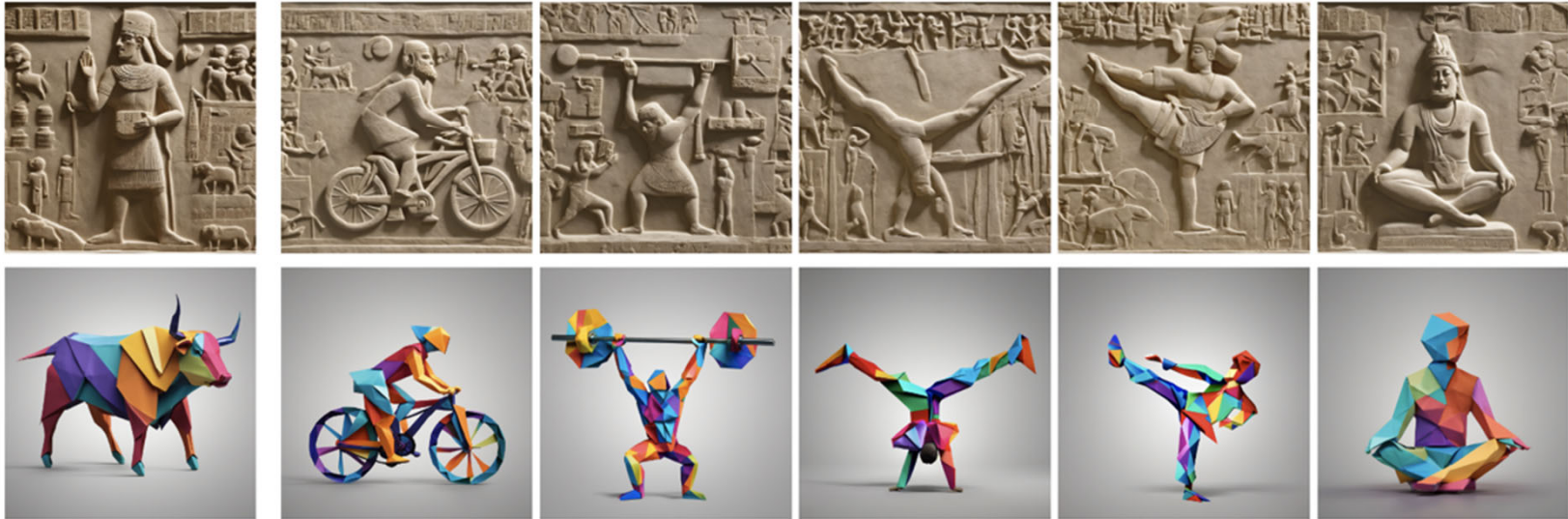


ControlNet + StyleAligned

Pose condition



Reference image



Textual Inversion+Dreambooth

Personalized content



"<V object> in the style of a beautiful paper-cut art."

+ StyleAligned

Personalized
content



Reference image



"<V object> in the style of a beautiful paper-cut art."

W.O AdalN

Personalized content



Reference image



"<V object> in the style of a beautiful paper-cut art."

DreamBooth + StyleAligned

Personalized content



Reference image



MultiDiffusion + StyleAligned

Reference image



"A poster in a flat design style."



"Houses in a flat design style."



"Mountains in a flat design style."



"Girrafes in a flat design style."

MultiDiffusion + StyleAligned

Reference image



"A poster in a flat design style."



W,O shared attention



W,O Attention Adaln



StyleAligned full

MultiDiffusion in Multi Styles

Left Reference



Right Reference



MultiDiffusion in Multi Styles

Left reference



Right reference





Left Reference

Right Reference



