

Images & Text

Rafail Fridman for DL4CV

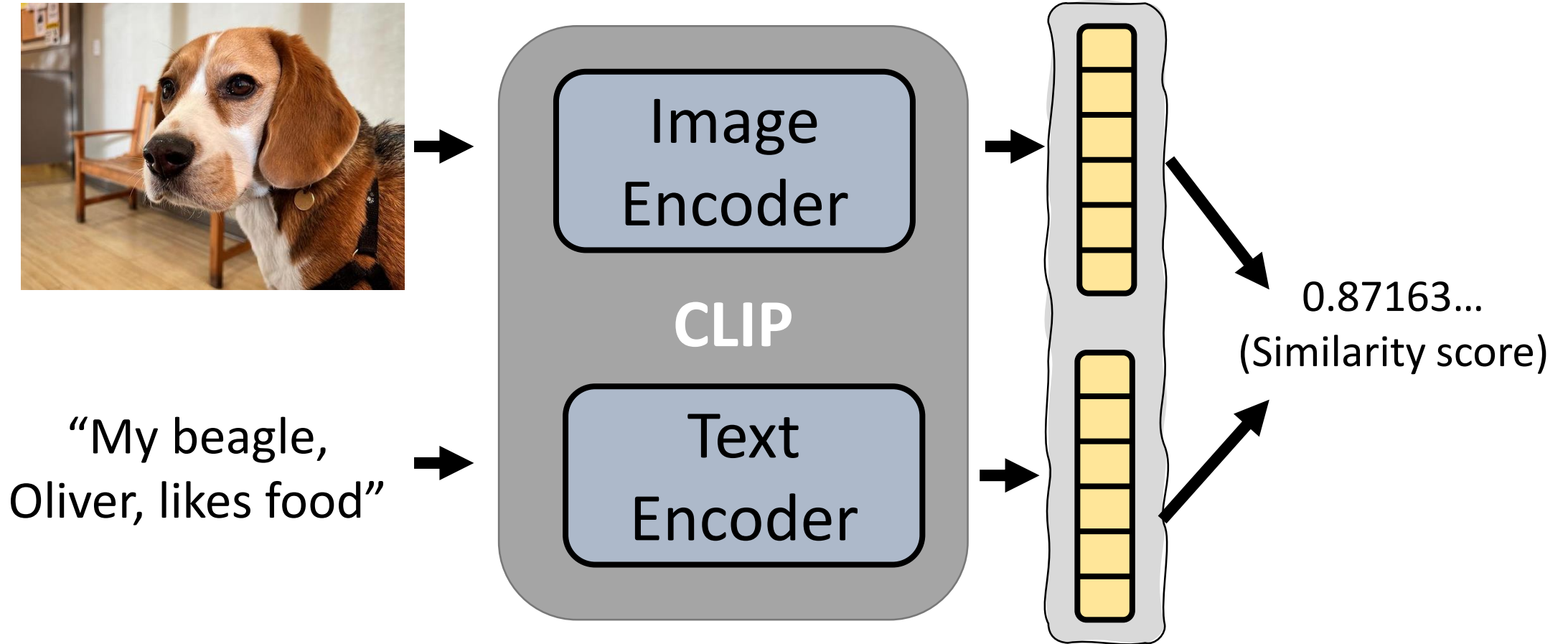
08.02.2023

Topics

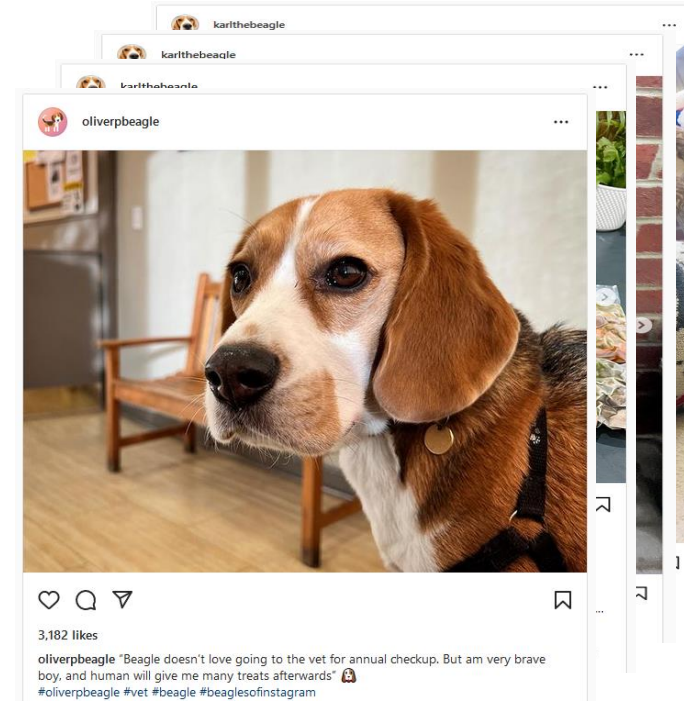
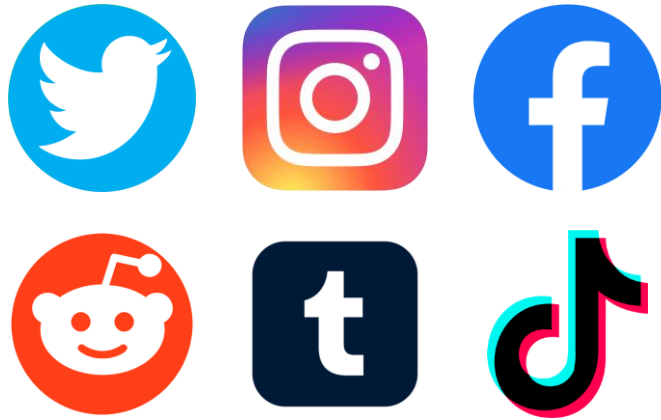
1. CLIP-guided optimization:
 - VQ-GAN + CLIP
 - StyleCLIP
 - Text2LIVE
2. Diffusion Models + text
 - Text conditioning in Diffusion Models
 - Classifier (free) guidance
 - Latent Diffusion models

CLIP - reminder

- Contrastive Language Image Pretraining



CLIP - reminder

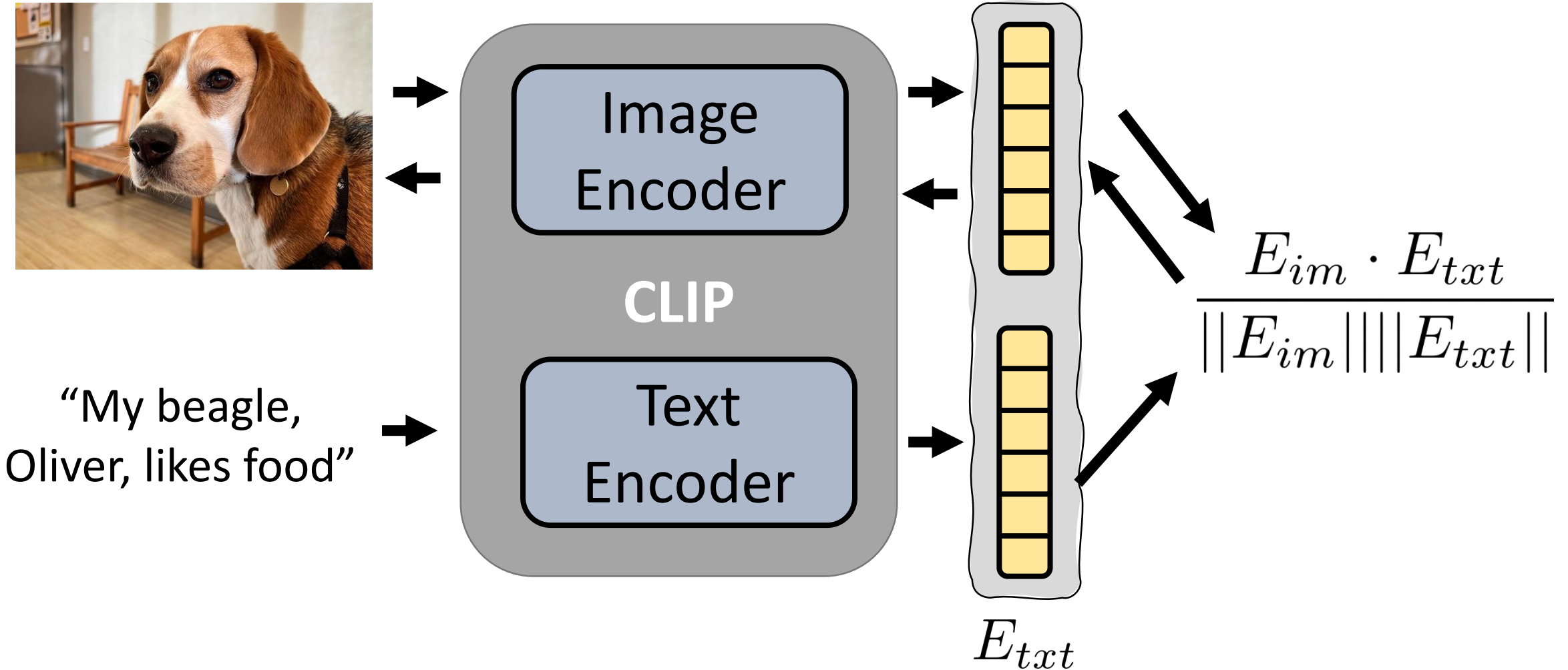


oliverpbeagle "Beagle doesn't love going to the vet for annual checkup. But am very brave boy, and human will give me many treats afterwards" 🐶 #oliverpbeagle #vet #beagle #beaglesofinstagram



**× 400
Million**

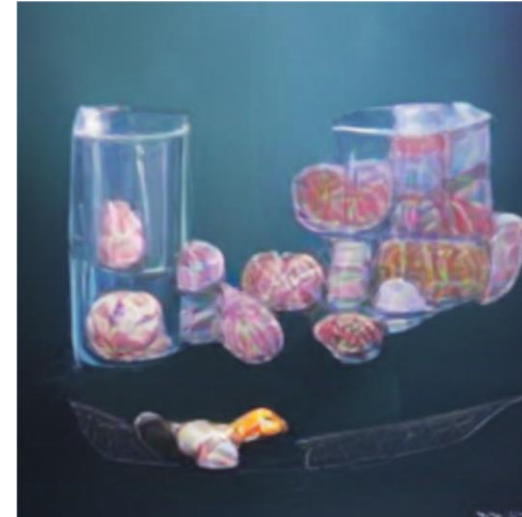
CLIP for generative tasks?



Using CLIP for generative tasks

Generation

A beautiful
painting of a
building in a
serene landscape



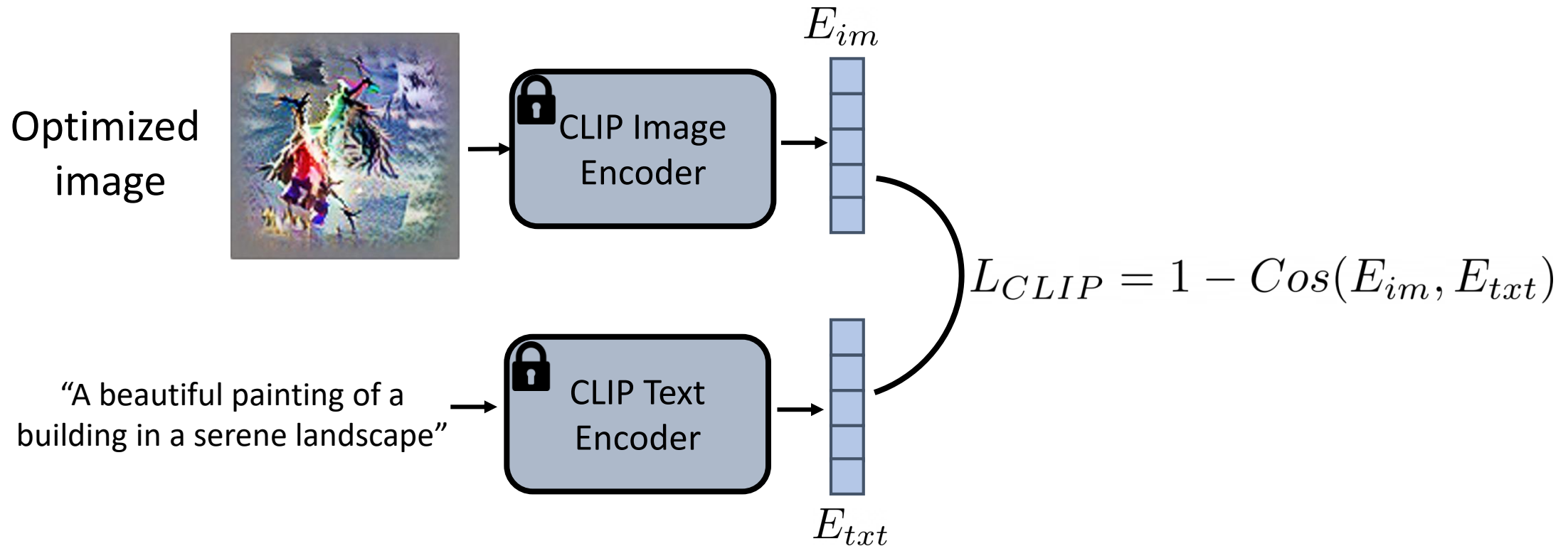
Editing



“A cake made of ice”

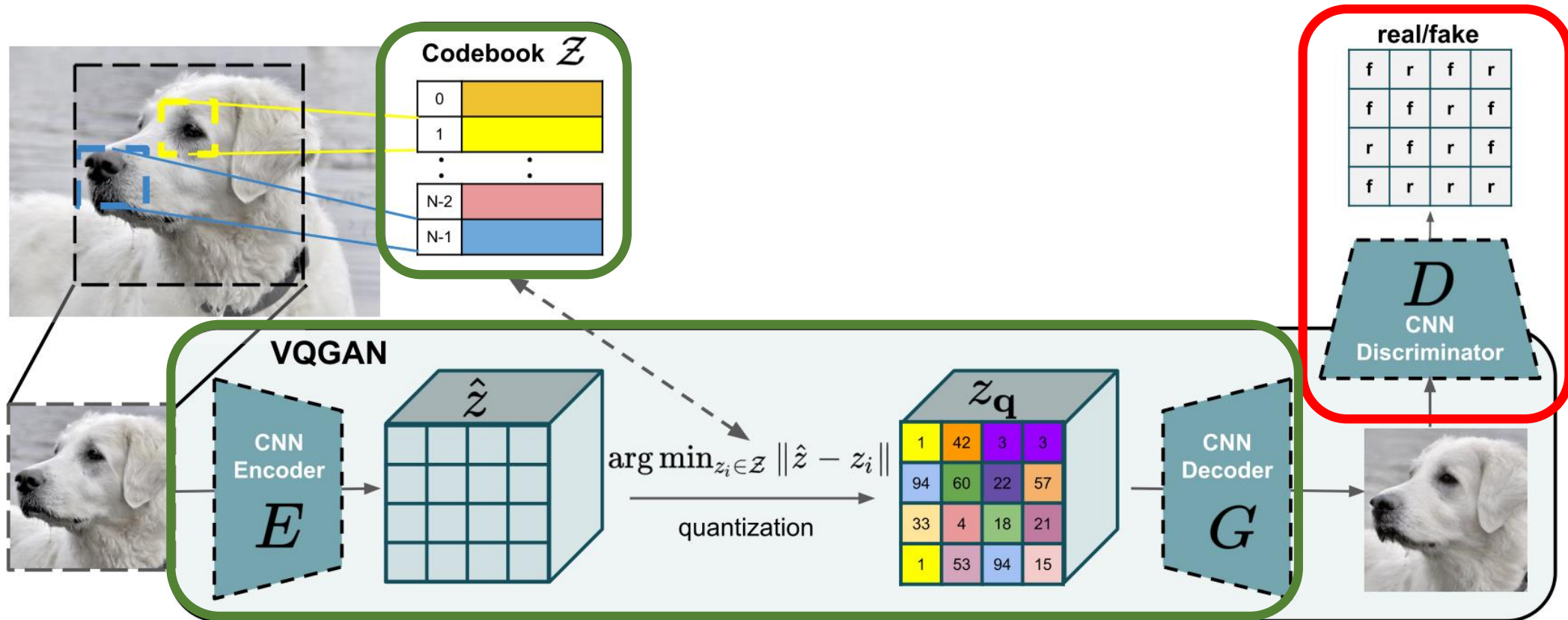


Naive approach



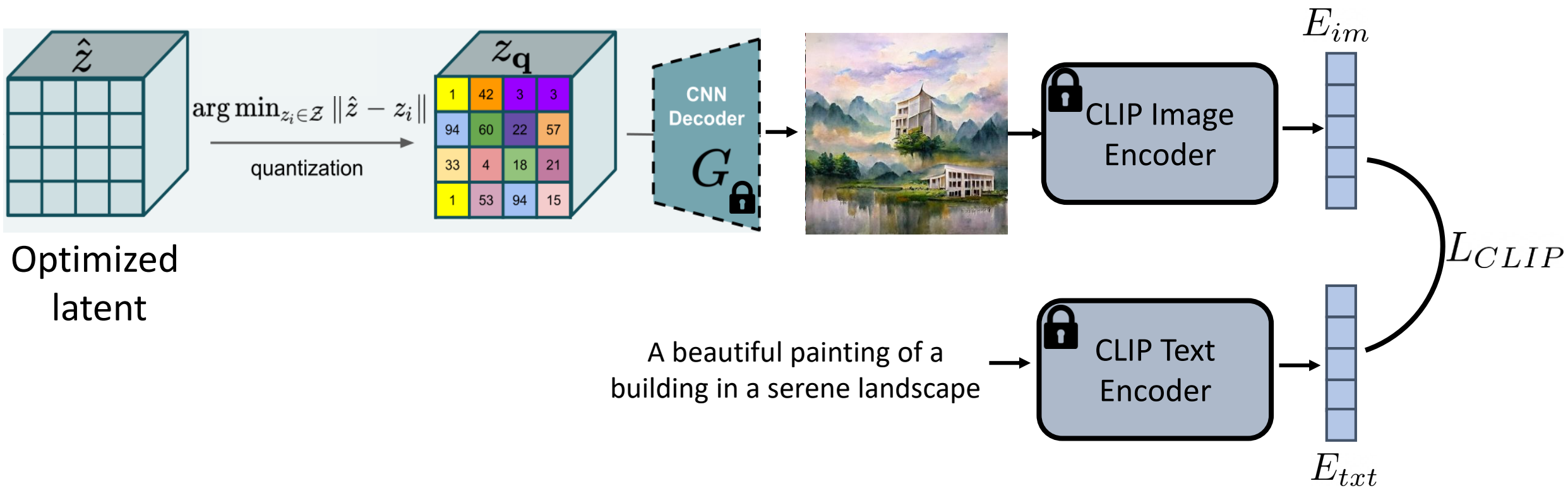
Utilize CLIP to steer the generation towards the desired text prompt

VQ-GAN - reminder



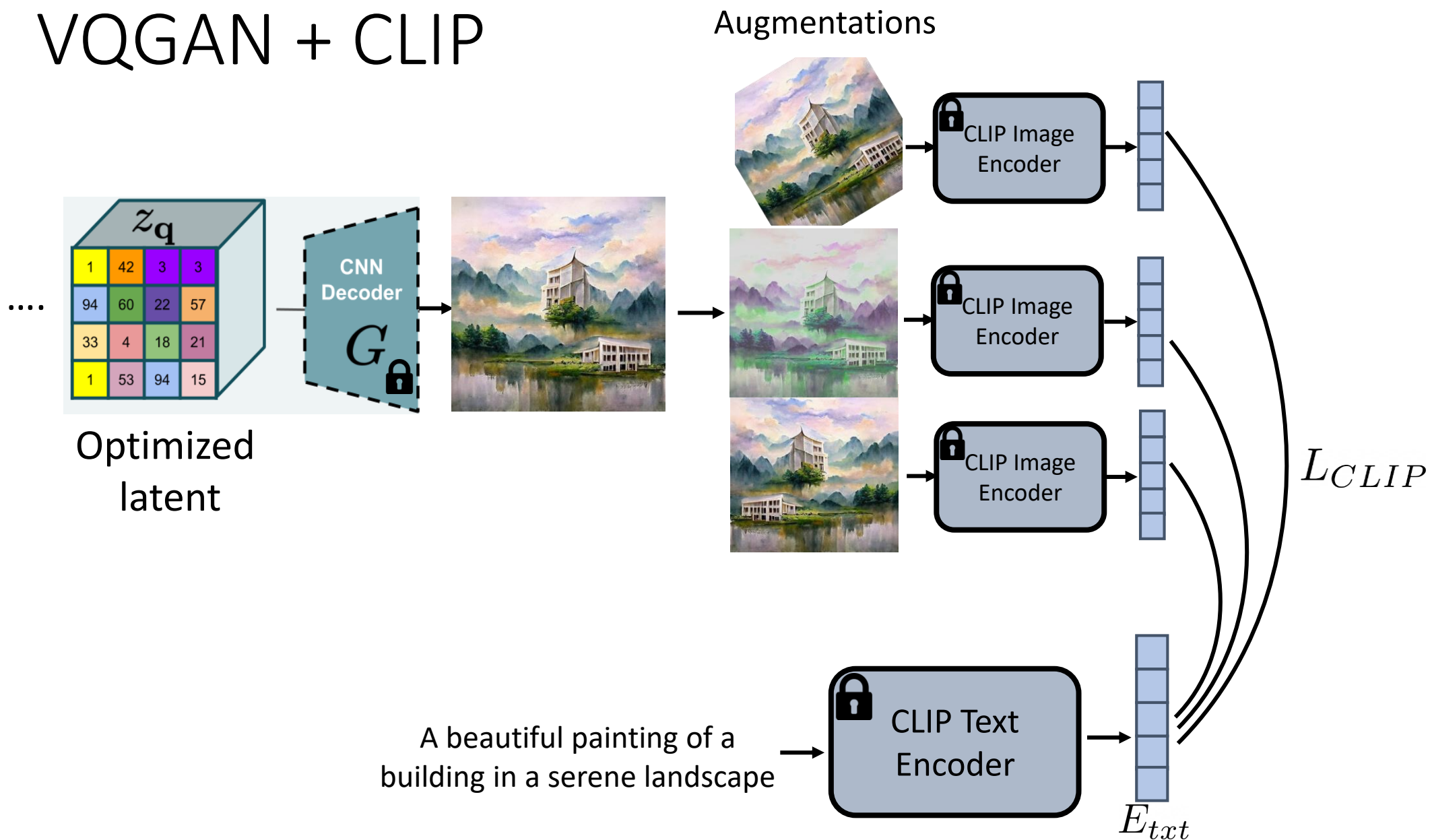
VQ-VAE + GAN

VQGAN + CLIP

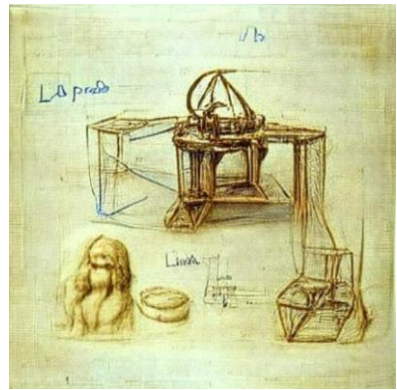


$$L_{CLIP} = 1 - \text{Cos}(E_{im}, E_{txt})$$

VQGAN + CLIP



VQGAN + CLIP results



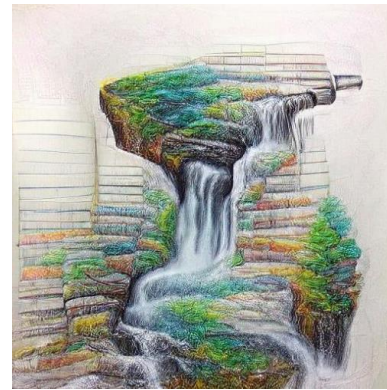
“A sketch of 3D printer by da Vinci”



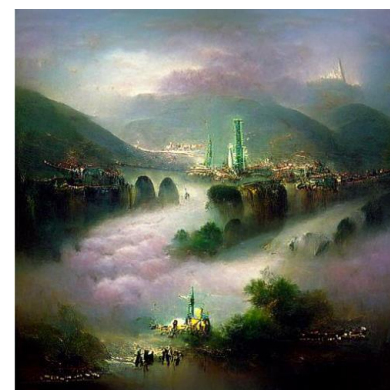
“An autogyro flying, artstation”



“A futuristic city in synthwave style”



“A colored pencil drawing of a waterfall”



“A painting of a city in a deep valley”

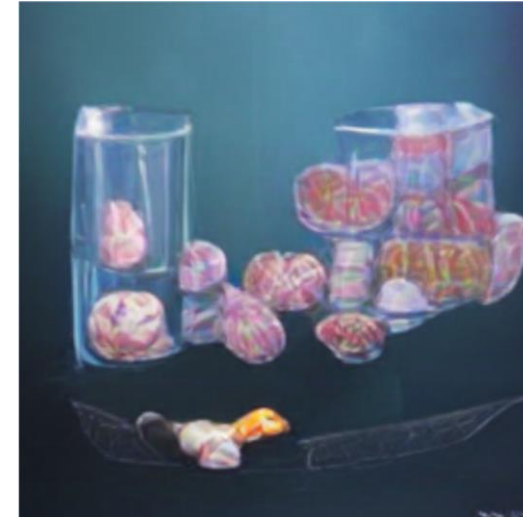


“Baba Yaga's house, fantasy art”

Using CLIP for generative tasks

Generation

A beautiful painting of a building in a serene landscape



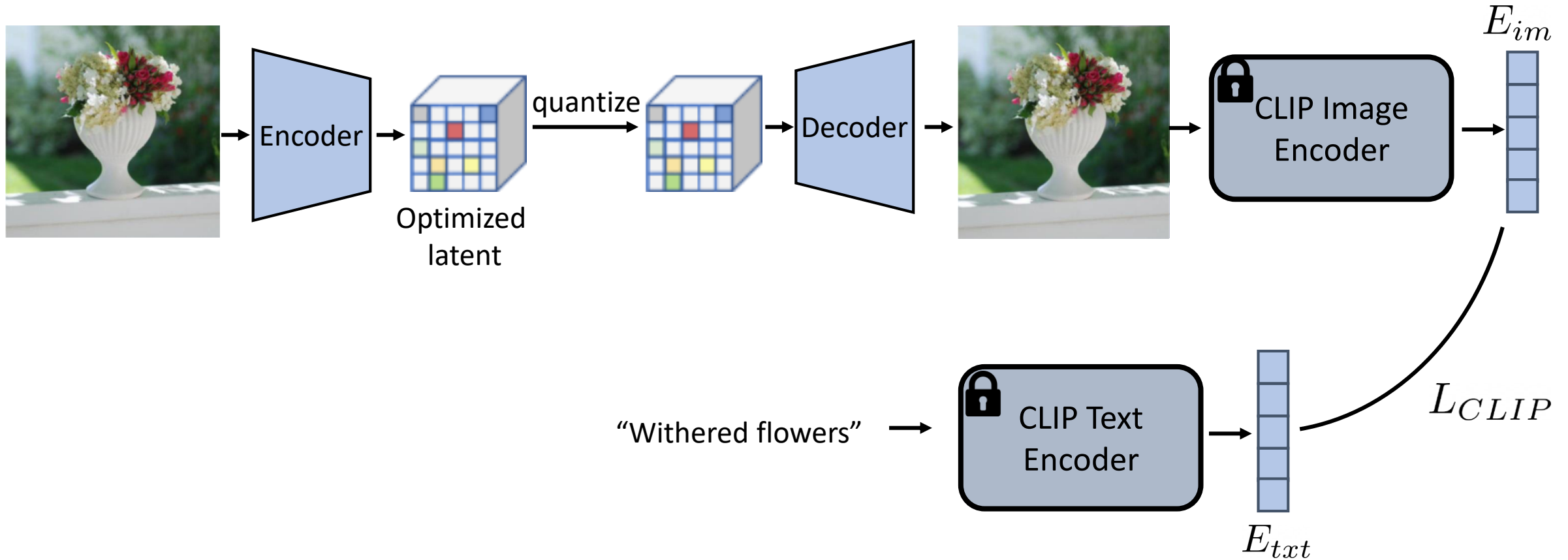
Editing



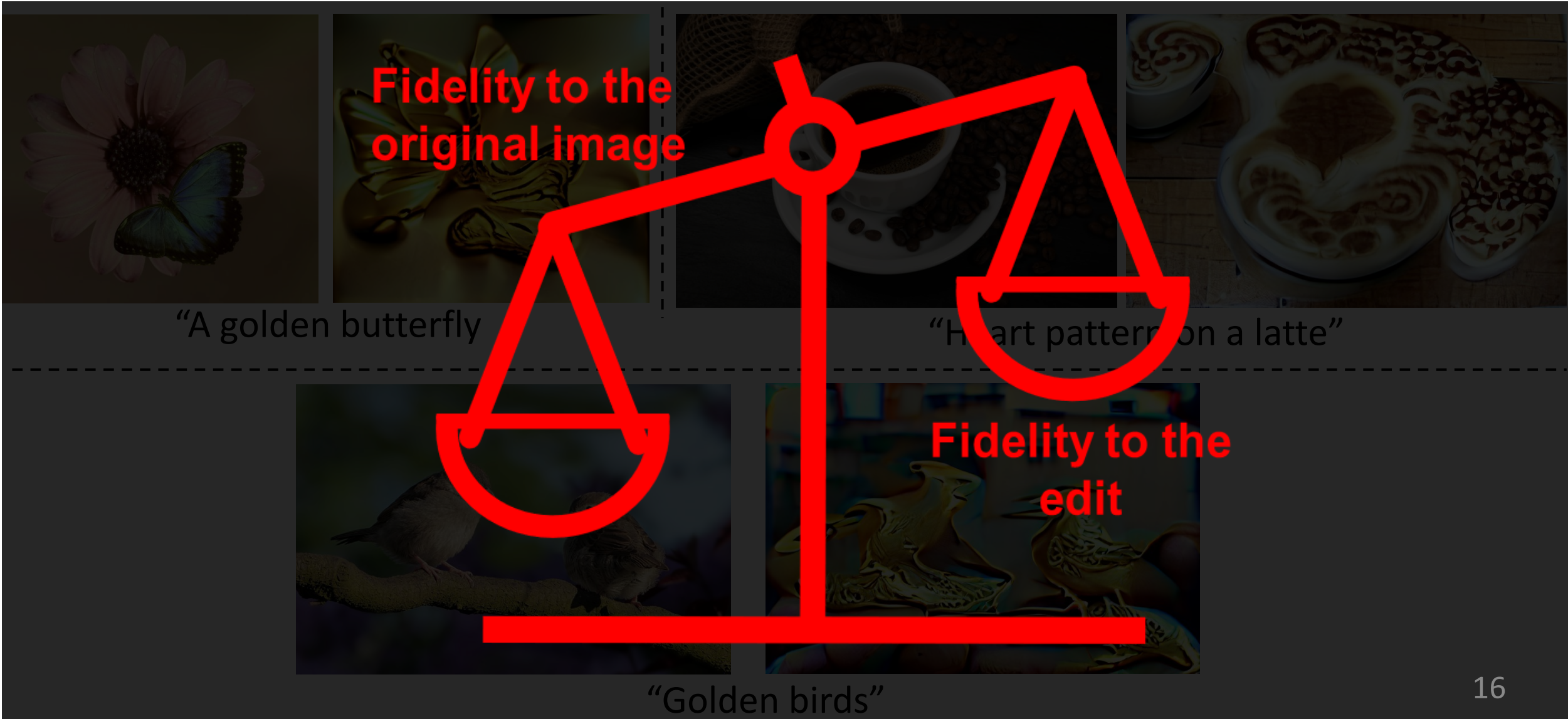
“A cake made of ice”



VQGAN + CLIP editing



VQGAN + CLIP editing



StyleCLIP – goal



“curly
hair”

“hi-top
fade hair”

“fringe
hair”

“black
hair”

“makeup”



“happy”

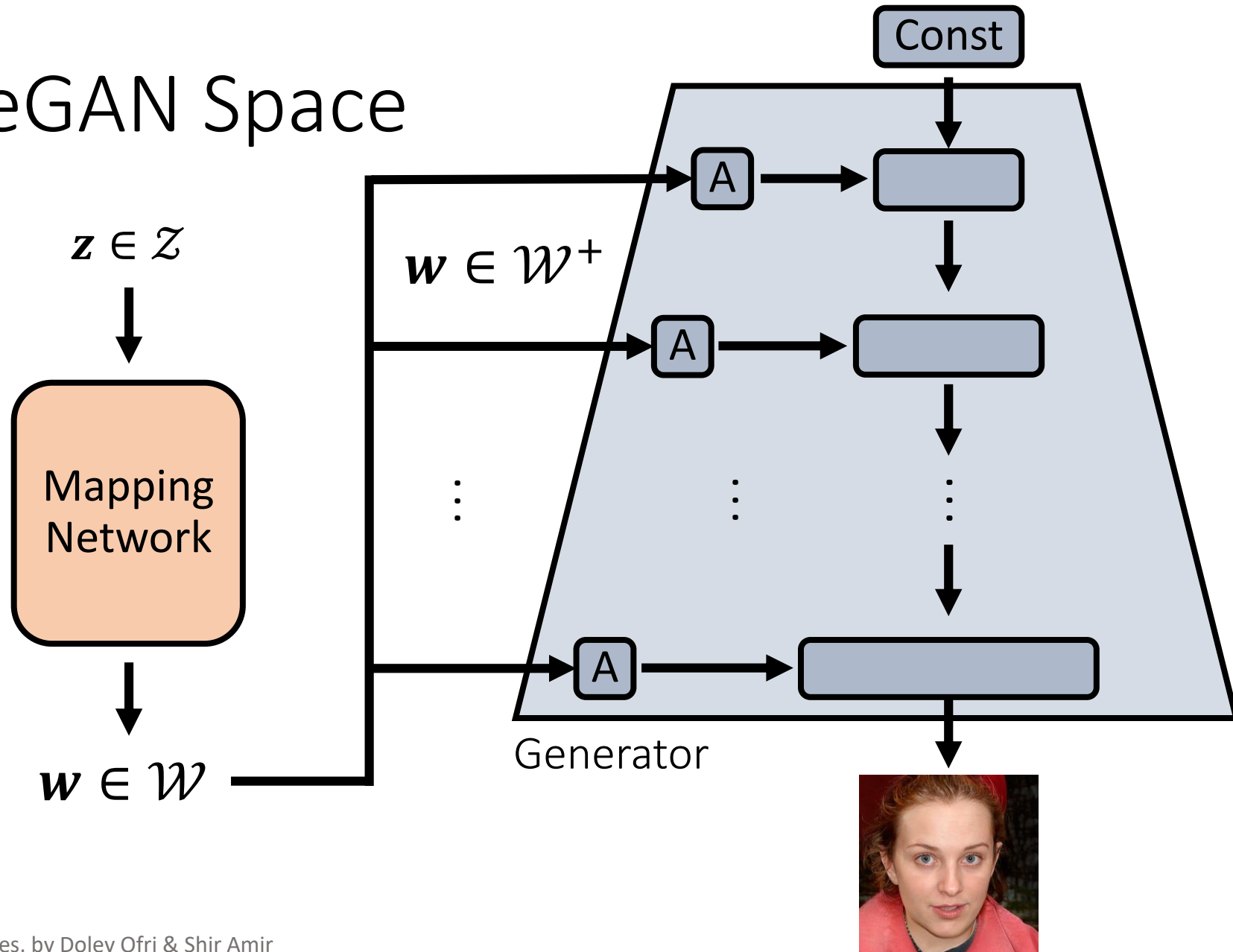
“big eyes”



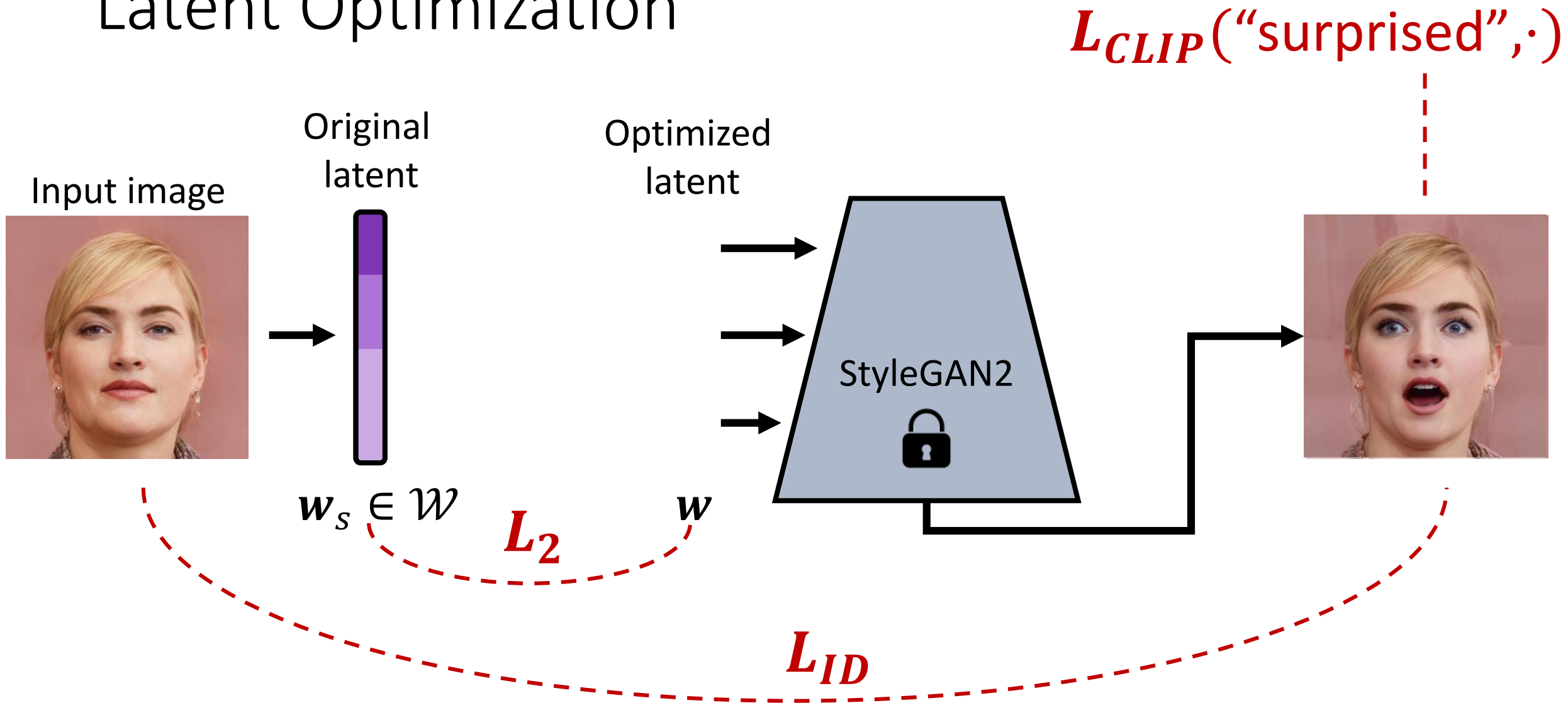
“cloud”

“spires”

StyleGAN Space

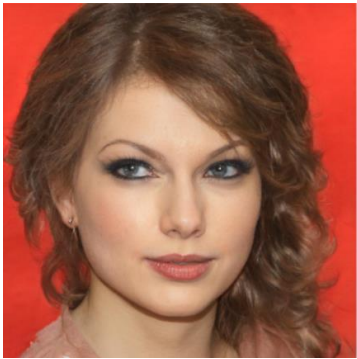


Latent Optimization

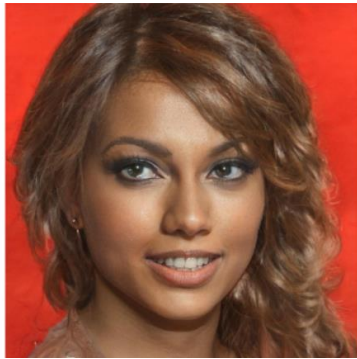


Adapted from ADLV2022 slides, by Dolev Ofri & Shir Amir

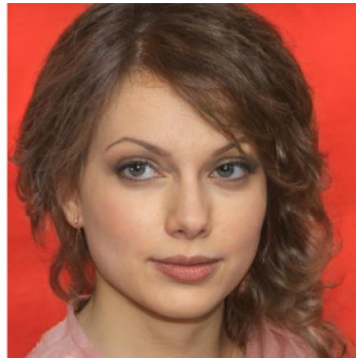
StyleCLIP - results



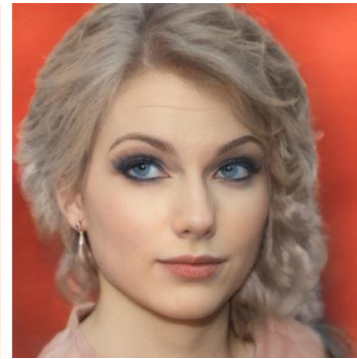
Input



"Beyonce"



"A woman
without
makeup"



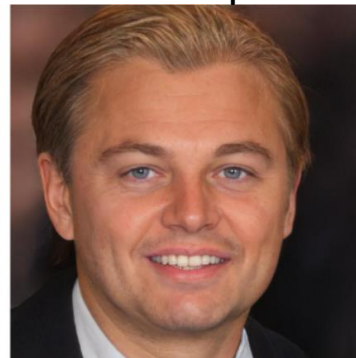
"Elsa from Frozen"



Input



"A man with
a beard"



"A blonde man"



"Donald Trump"

CLIP knows about these!

StyleCLIP – results, different domains



StyleCLIP – results, different domains

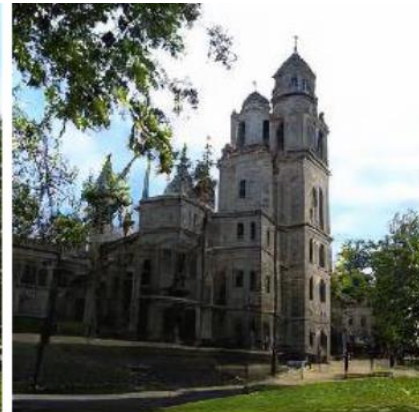
Input

Trees

Clouds

Spires

Round Roof



Text2LIVE – goal

Original image



“oreo cake”



“brioche”



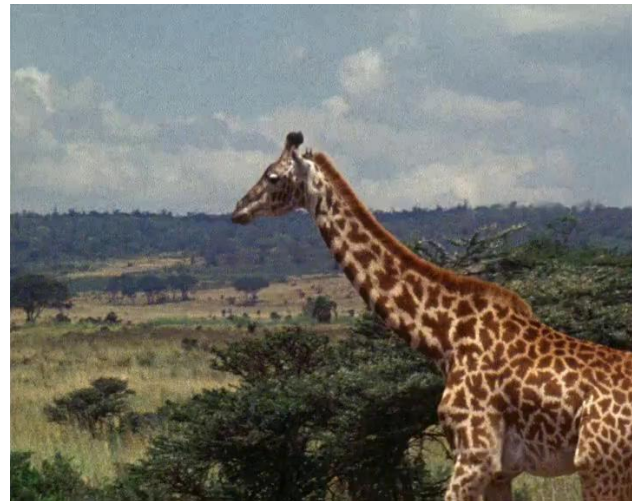
“ice”



“spinach moss cake”



Original video



“stained glass giraffe”



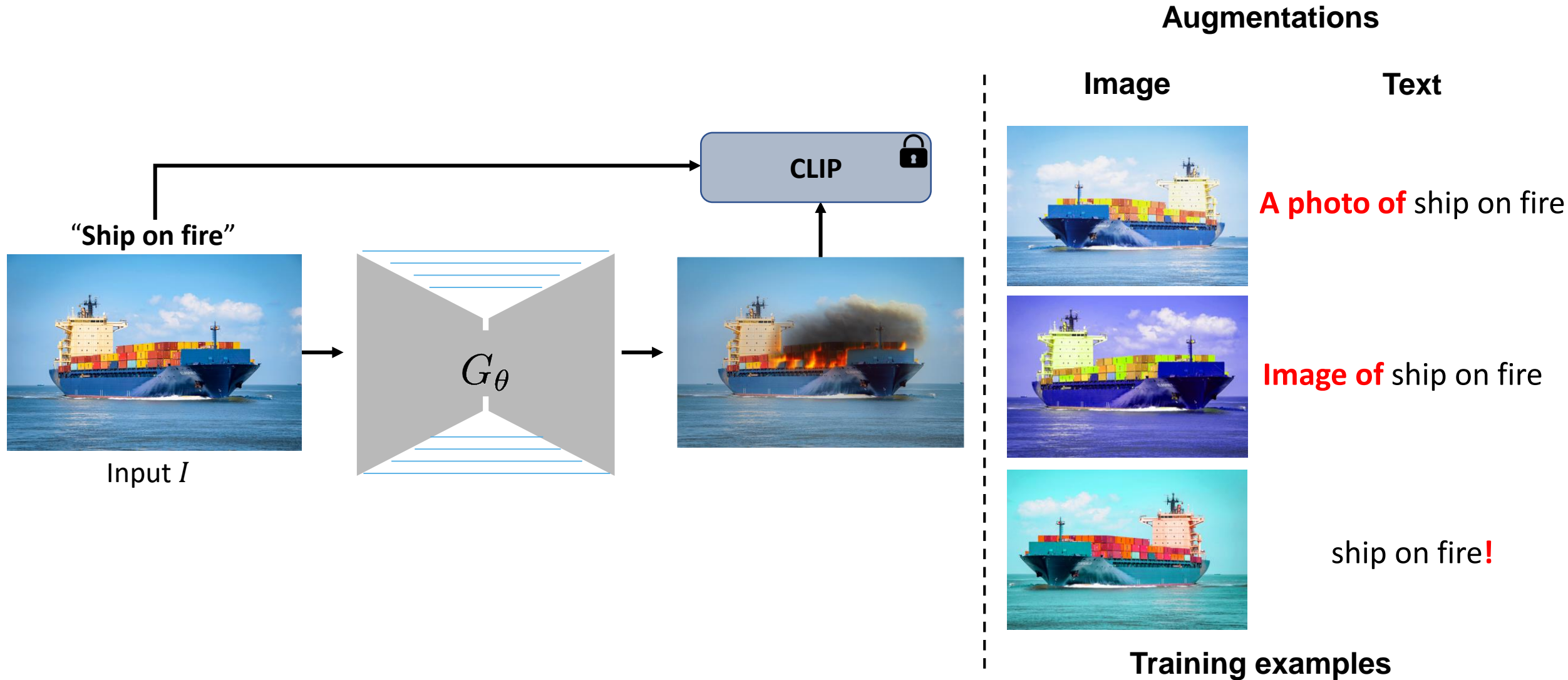
“giraffe with neck warmer”



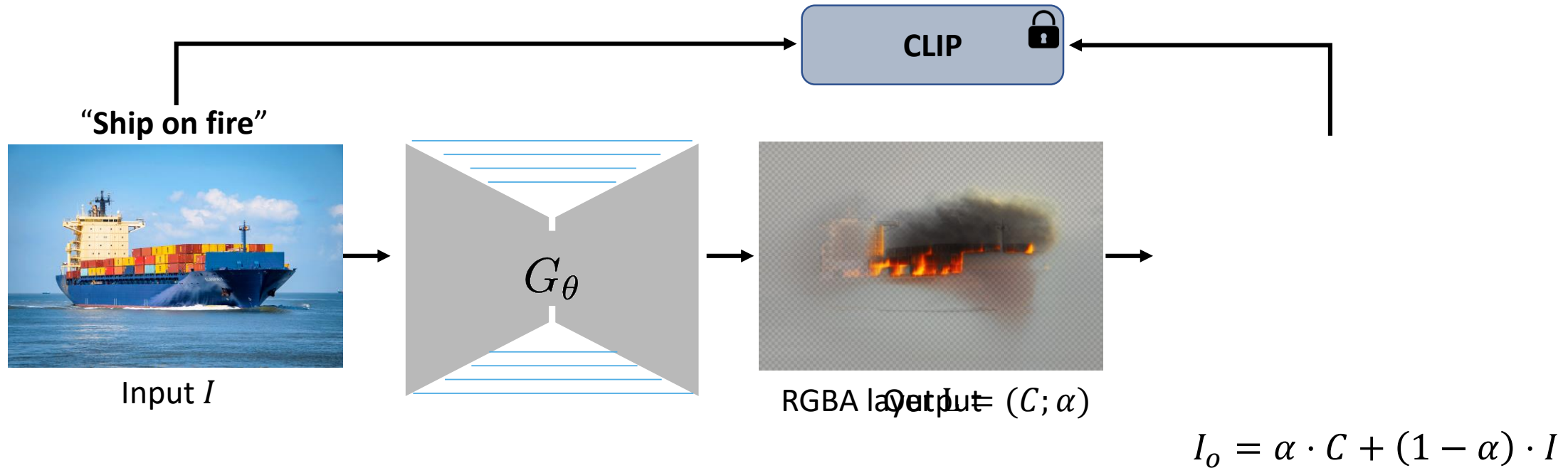
“giraffe with hairy colorful mane”



Text2LIVE – method



Text2LIVE – method

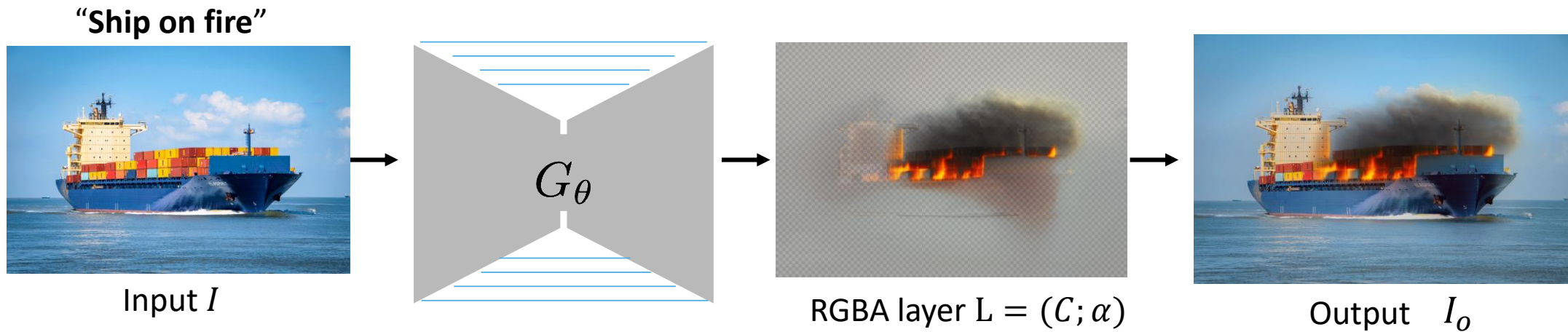


Text2LIVE – losses

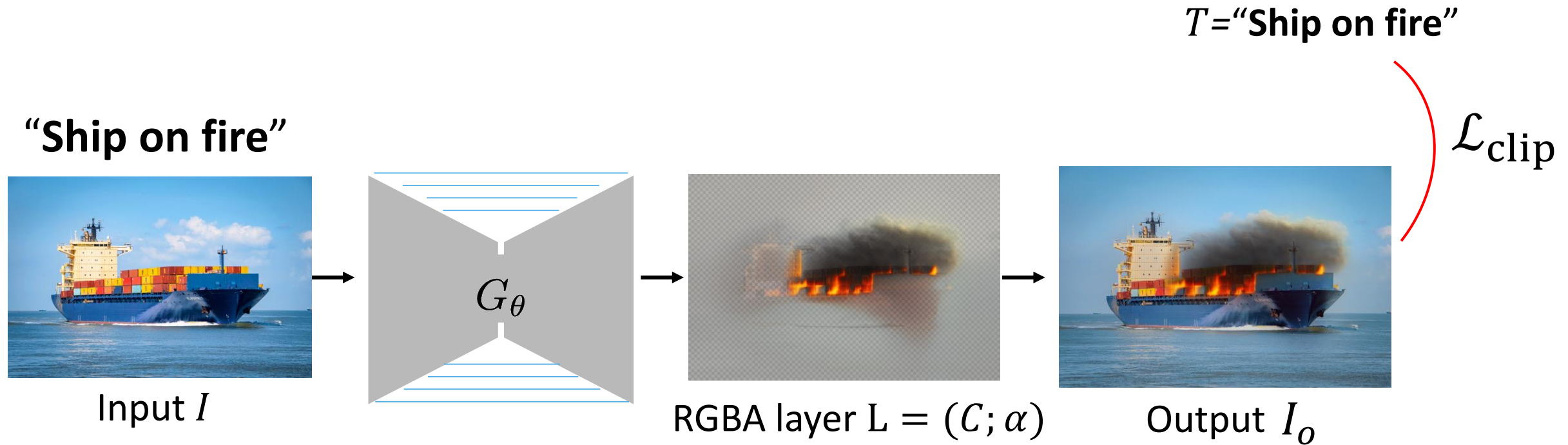
$$\mathcal{L}_{Text2LIVE} = \underbrace{\mathcal{L}_{comp}(I_o)}_{\text{Final composite}} + \underbrace{\alpha \mathcal{L}_{screen}(C, \alpha) + \gamma \mathcal{L}_{sparsity}(\alpha)}_{\text{Edit layer}}$$

Final
composite

Edit layer



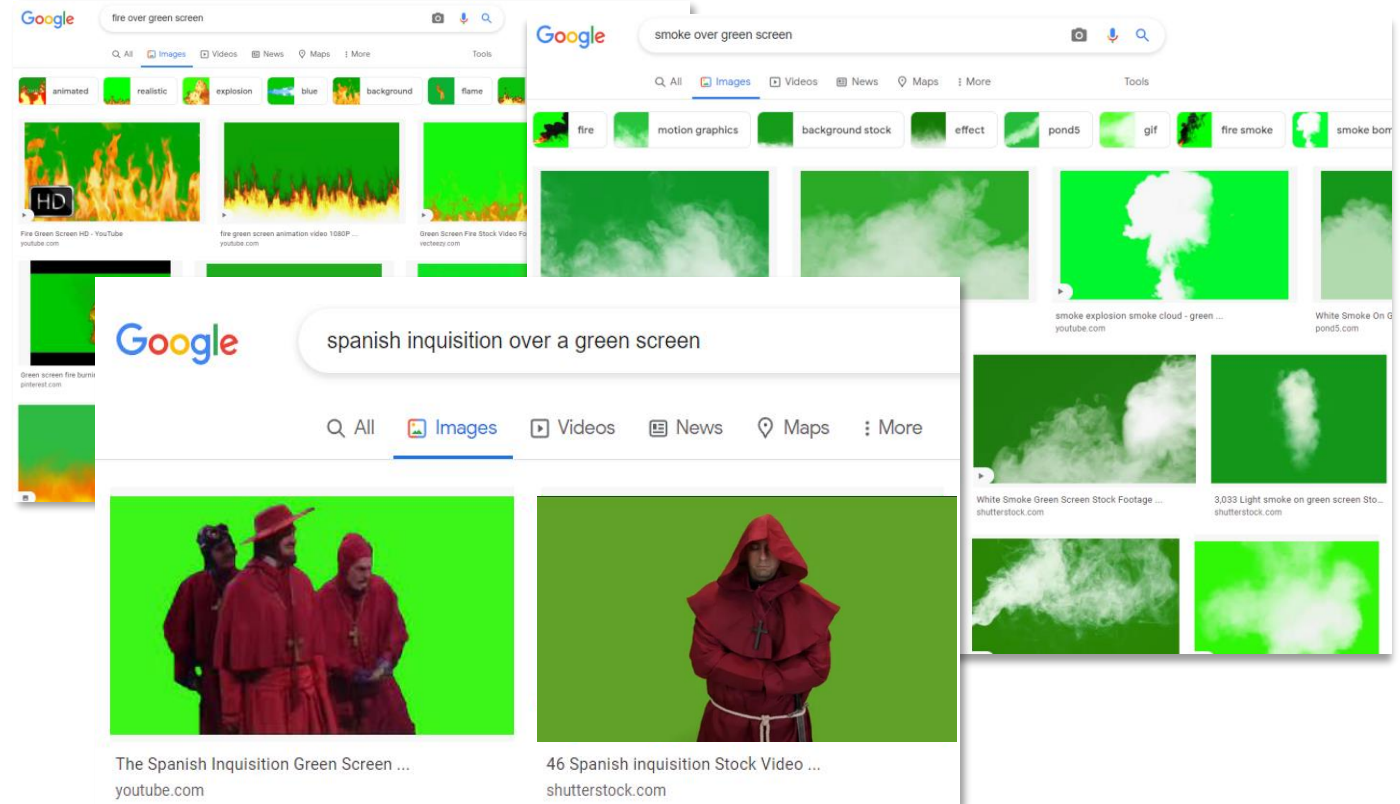
Loss on the composition



Losses on the edit layer

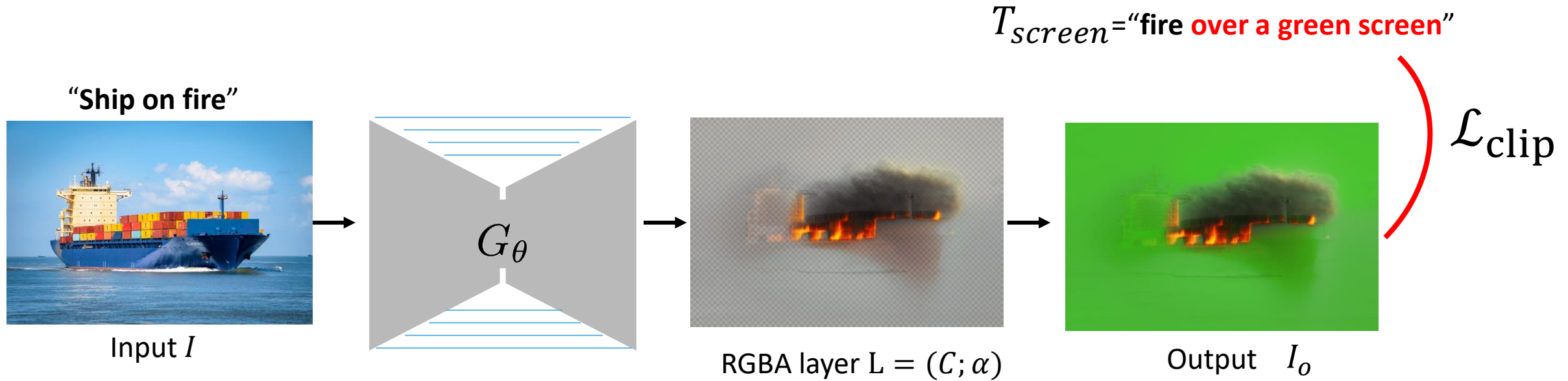
$$\mathcal{L}_{Text2LIVE} = \mathcal{L}_{comp}(I_o) + \alpha \mathcal{L}_{screen}(C, \alpha) + \gamma \mathcal{L}_{sparsity}(\alpha)$$

Chroma keying



Text2LIVE – losses

$$\mathcal{L}_{Text2LIVE} = \mathcal{L}_{comp}(I_o) + \alpha \mathcal{L}_{screen}(C, \alpha) + \gamma \mathcal{L}_{sparsity}(\alpha)$$



Losses on the edit layer

w/ \mathcal{L}_{screen} ($T = \text{"smoke"}$)

"A man smoking a cigar"



w/o \mathcal{L}_{screen}



Losses on the edit layer



Input Image

“woman wearing a red hat”



Relevancy map*
“hat”



Text2LIVE output matte

MSE



Text2LIVE result

Text2LIVE – results

Input



“wooden *”



“golden *”



“stained glass *”

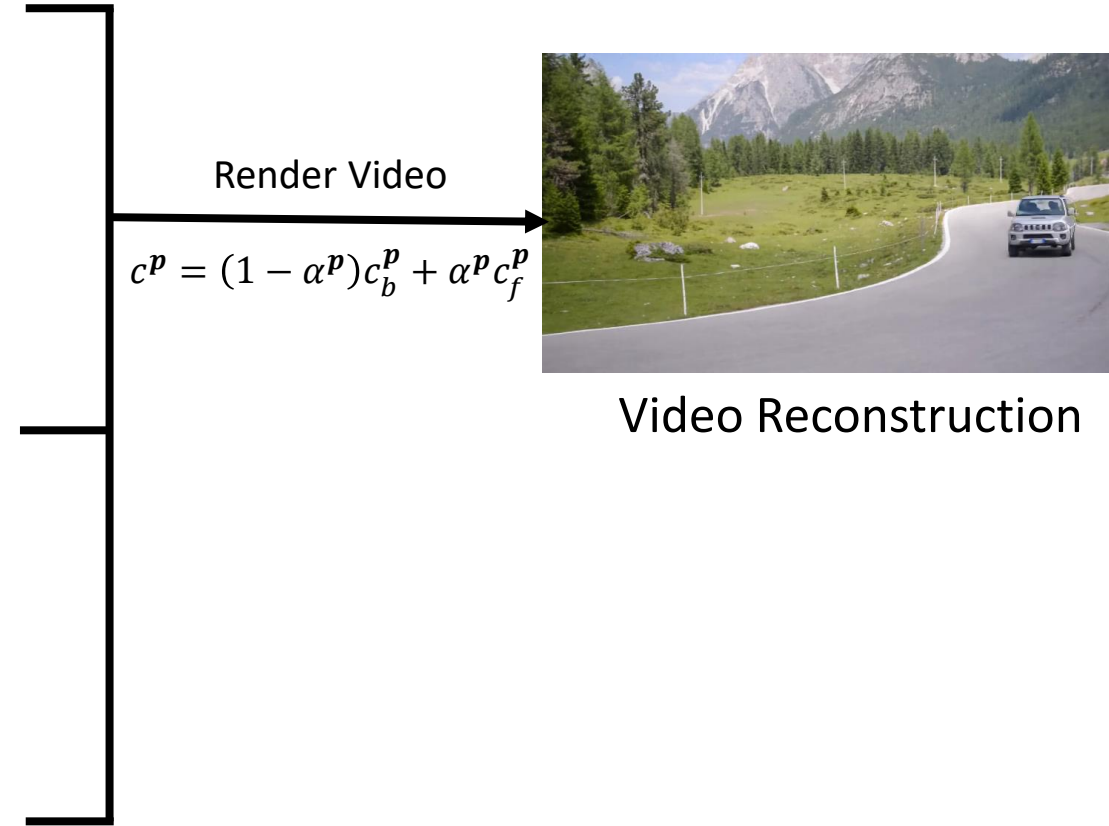


“crochet *”



NLA - reminder

Input video



Video Reconstruction

NLA - reminder

Input video



Render Video

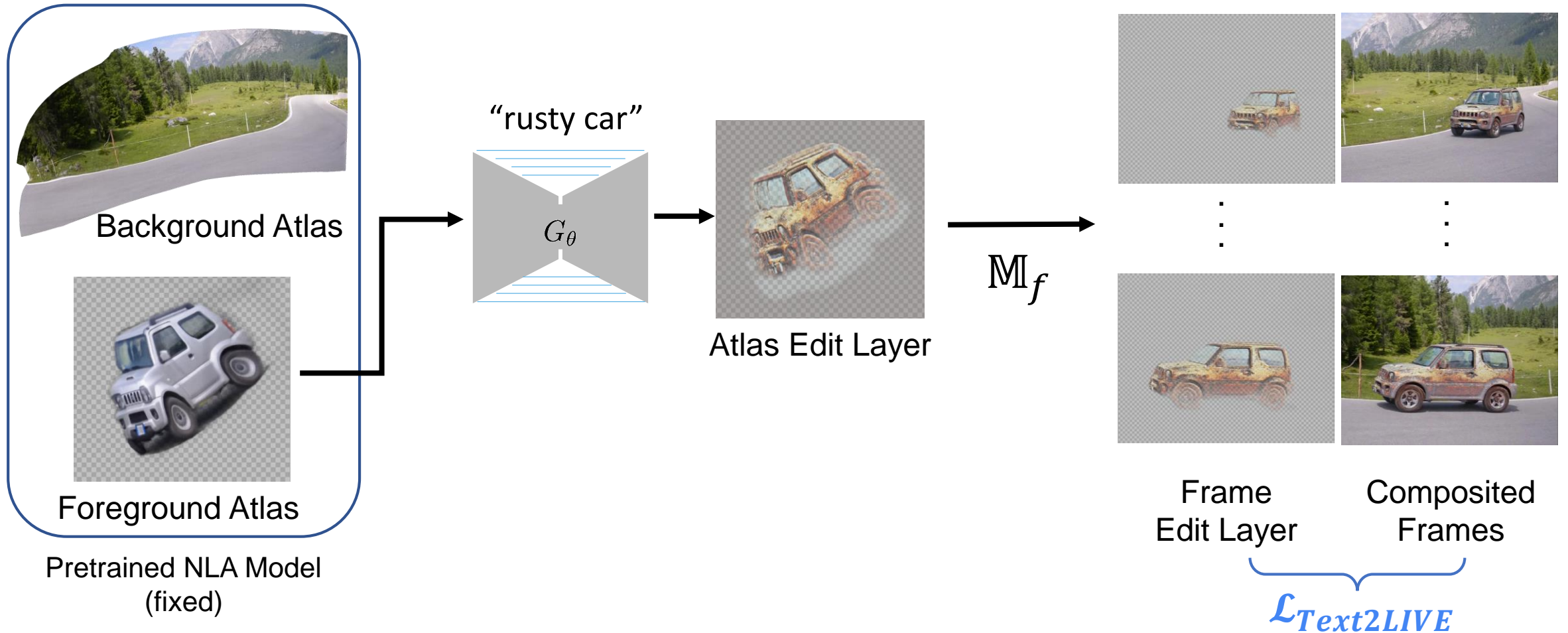
$$c^p = (1 - \alpha^p)c_b^p + \alpha^p c_f^p$$



Video Reconstruction

Editing huge pixel volume → Editing a single 2D image

Text2LIVE video editing



Text2LIVE – results

Input Video



“swarovski blue crystal swan”



Original Video



Input Video



“dalmatian dog”



“dog with leopard texture”



Topics

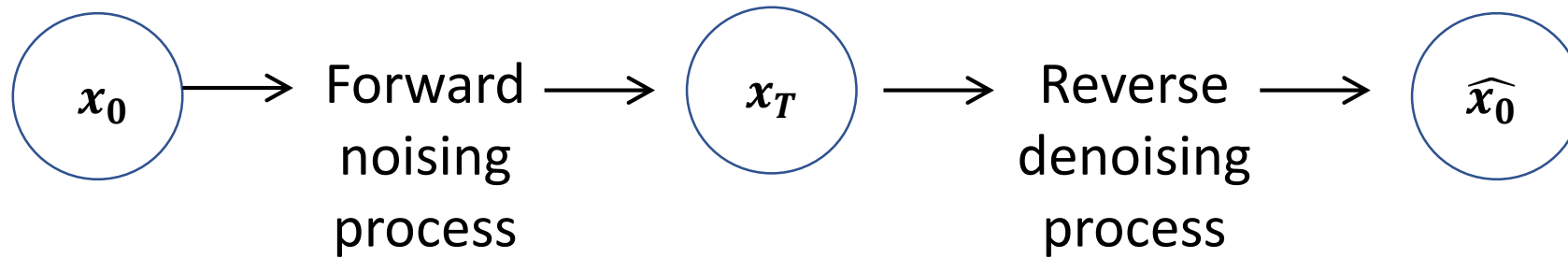
1. CLIP-guided optimization:

- VQ-GAN + CLIP
- StyleCLIP
- Text2LIVE

2. **Diffusion Models + text**

- **Text conditioning in Diffusion Models**
- **Classifier (free) guidance**
- **Latent Diffusion models**

Diffusion Models - reminder



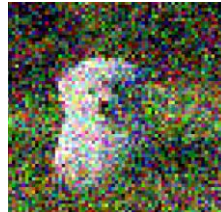
Reverse
denoising
process



x_0



x_1




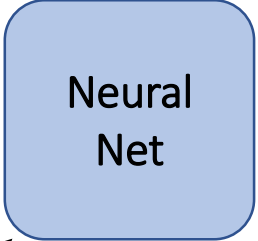
x_T

Generating samples by gradually reducing noise

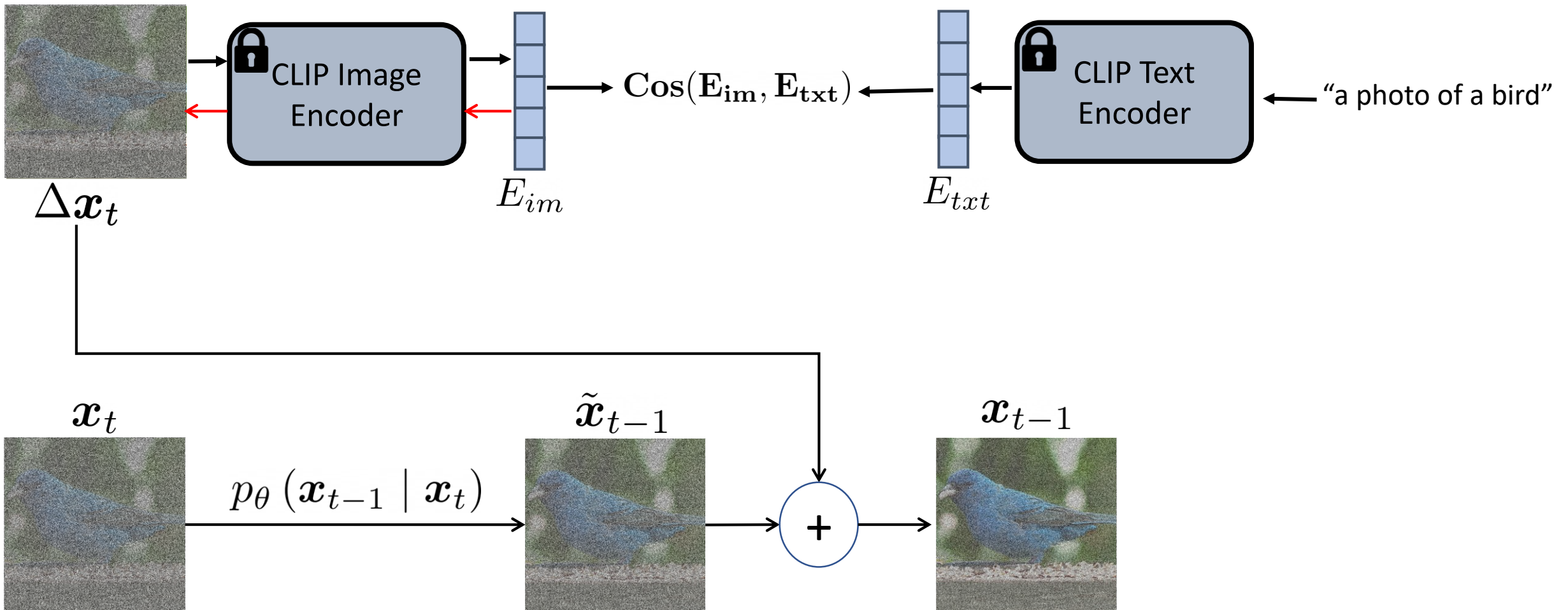
Reminder:

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \sigma_t^2 \mathbf{I})$$

Sampling Algorithm

1. Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2. For $t = T, \dots, 1$:
 - “Predict noise” in current image $\mathbf{z}_{\theta}(\mathbf{x}_t, t)$ 
 - Sample variance noise $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - Sample from reverse distribution:
$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{z}_{\theta}(\mathbf{x}_t, t) \right) + \tilde{\beta}_t \mathbf{z}$$
3. Return \mathbf{x}_0

Conditional image generation **with CLIP**



Conditional sampling with CLIP model

1. Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, text prompt c , text encoder f , image encoder g

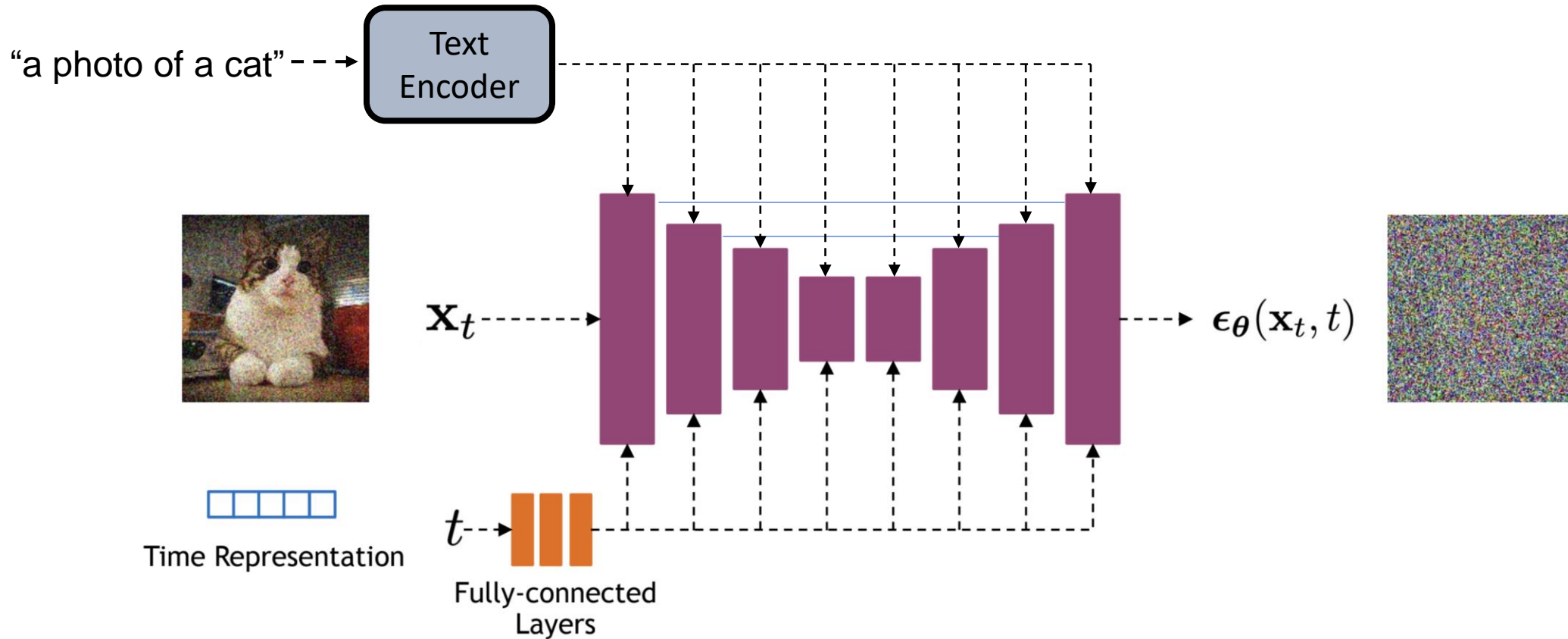
2. For $t = T, \dots, 1$:

- “Predict noise” in current image $\mathbf{z}_\theta(\mathbf{x}_t, t)$
- Sample variance noise $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Sample from reverse distribution:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{z}_\theta(\mathbf{x}_t, t) \right) + \tilde{\beta}_t \mathbf{z} + s \nabla_{x_t} (f(x_t) \cdot g(c))$$

3. Return \mathbf{x}_0

Text conditioning implementation

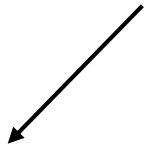


Time features are fed to the residual blocks using either simple spatial addition or using adaptive group normalization layers. (see Dhariwal and Nichol NeurIPS 2021)

Text conditioning in Diffusion Models



Noisy image x_t



$\phi(x_t)$

“A photo of a cat”

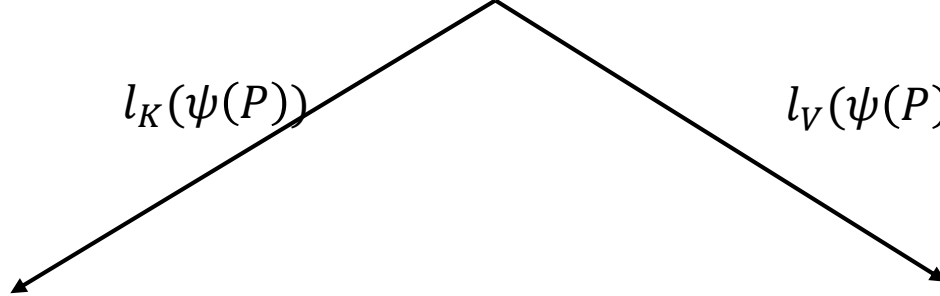
Prompt P



Text tokens $\psi(P)$

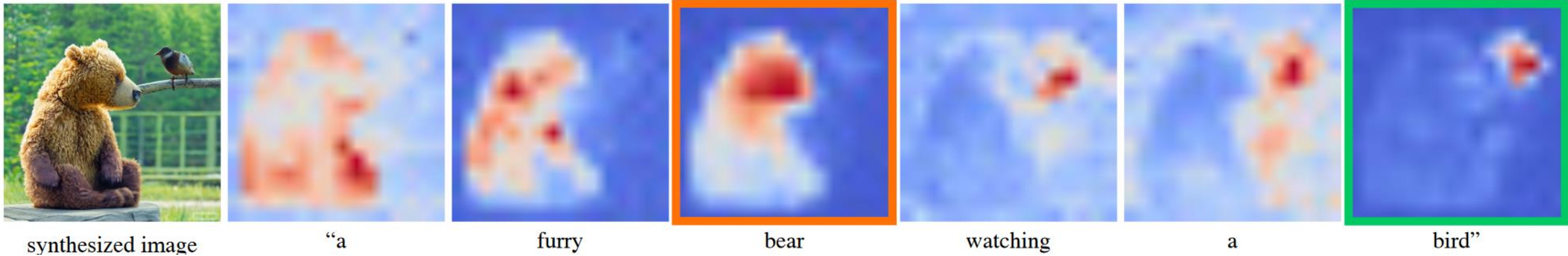
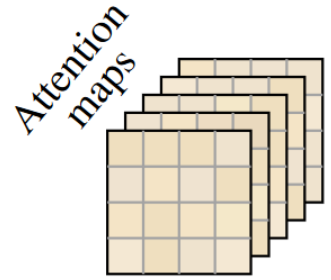
$l_K(\psi(P))$

$l_V(\psi(P))$



$\hat{\phi}(x_t)$

Text conditioning in Diffusion Models



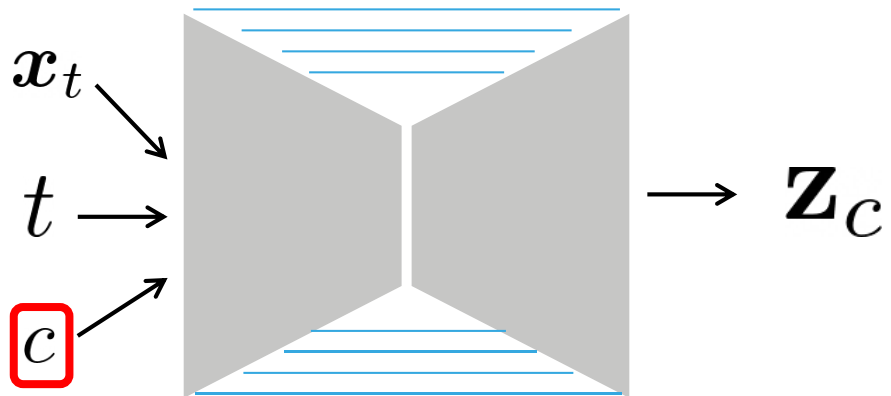
Average attention maps across all timestamps

Visual features attend to the text prompt tokens

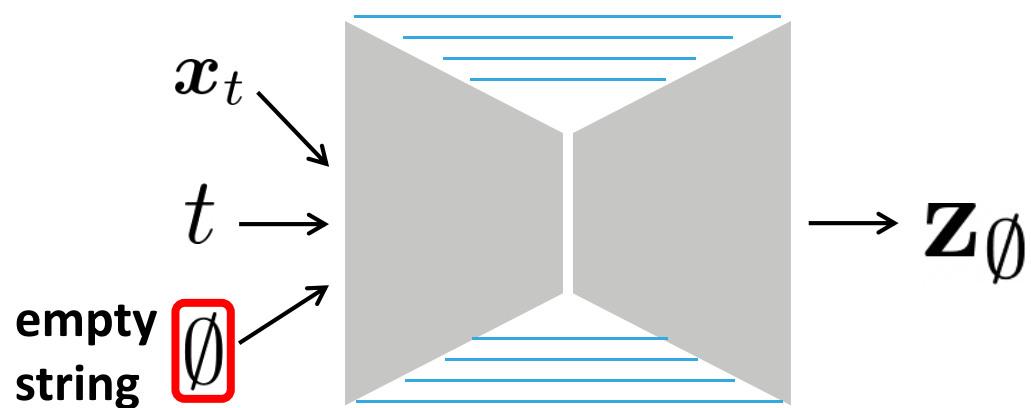
Classifier-free guidance

Training

Text conditioning

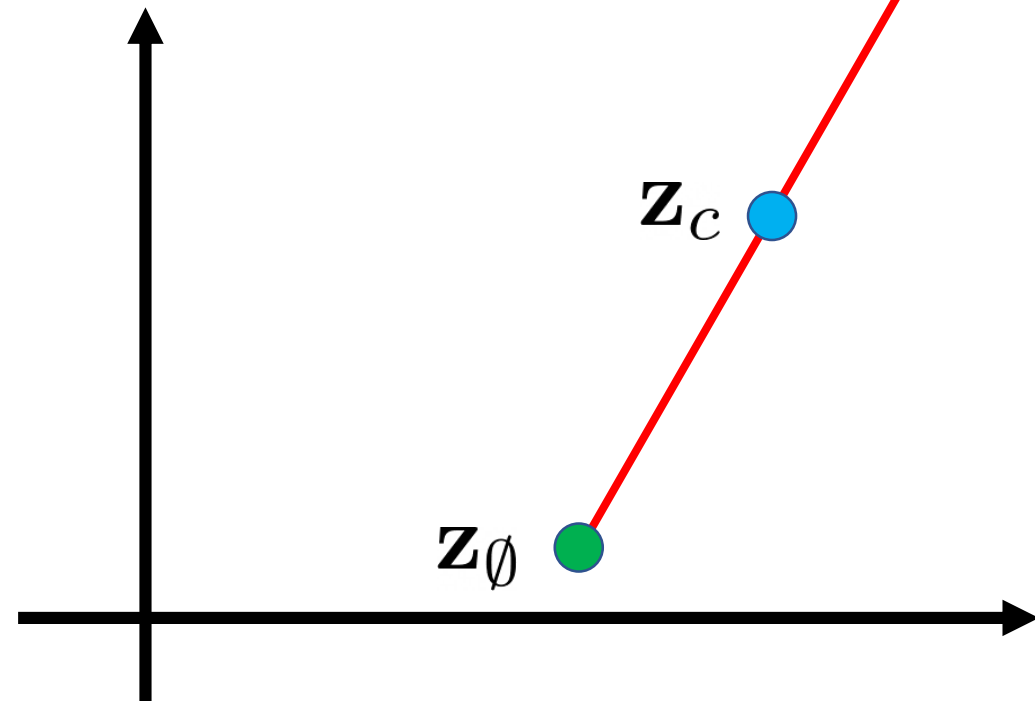


No conditioning



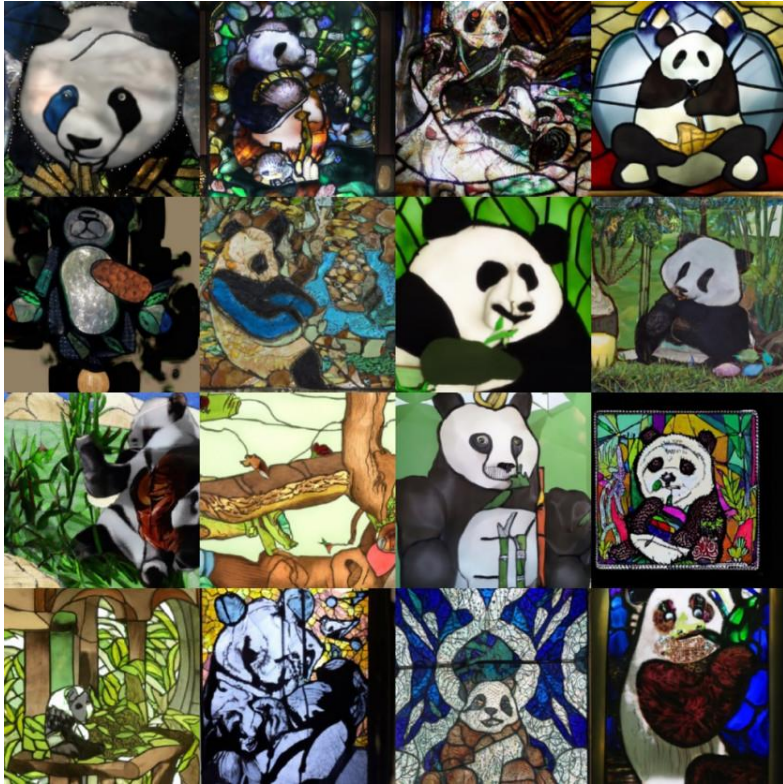
Inference

$$\mathbf{z}_\emptyset + s \cdot (\mathbf{z}_c - \mathbf{z}_\emptyset)$$



Effect of Classifier-free guidance

"A stained glass window of a panda eating bamboo"



$s = 1$

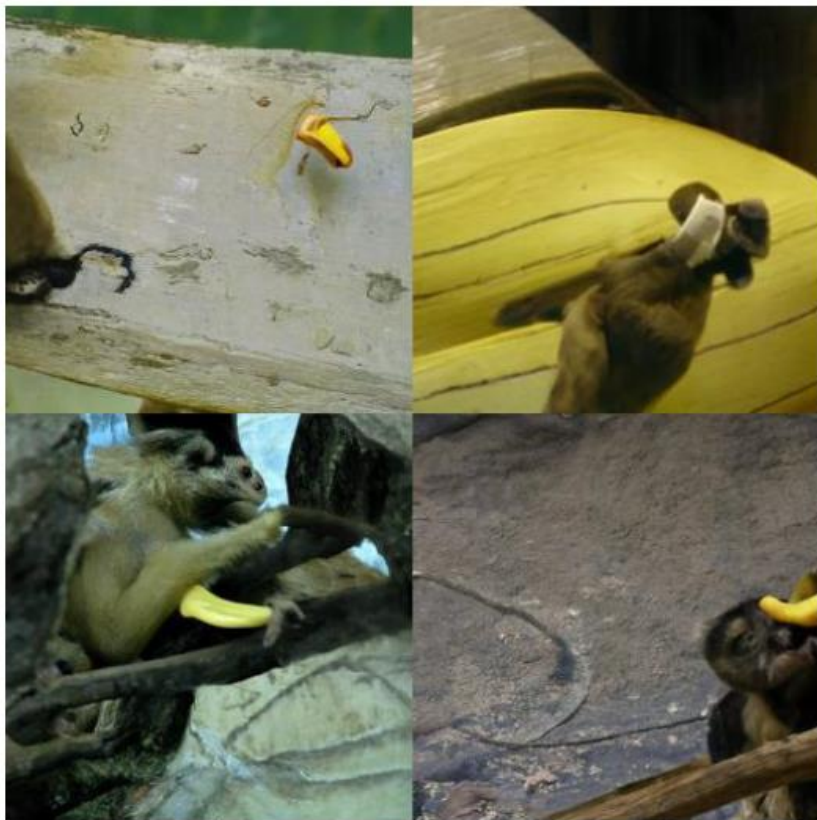


$s = 3$

$$\mathbf{z}_\emptyset + s \cdot (\mathbf{z}_c - \mathbf{z}_\emptyset)$$

Guidance comparison

A monkey eating a banana



CLIP guidance



Noised CLIP guidance



Classifier-free guidance

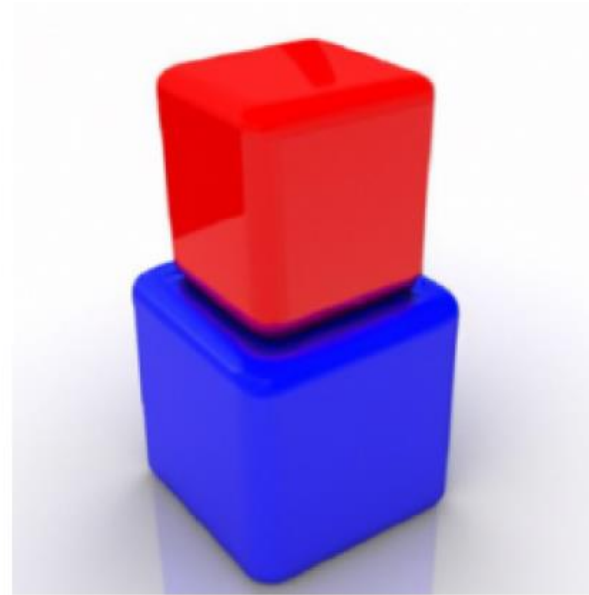
Results



“a hedgehog using a calculator”



“a corgi wearing a red bowtie and purple party hat”



“a red cube on top of a blue cube”



“a high-quality oil painting of a psychedelic hamster dragon”

One small problem

- Diffusion models operate in pixel space
- Training/inference on high-resolution images:
 - Long training time
 - A lot of GPU memory

Find a different space for diffusion models 😊



256 x 256

OK



512 x 512

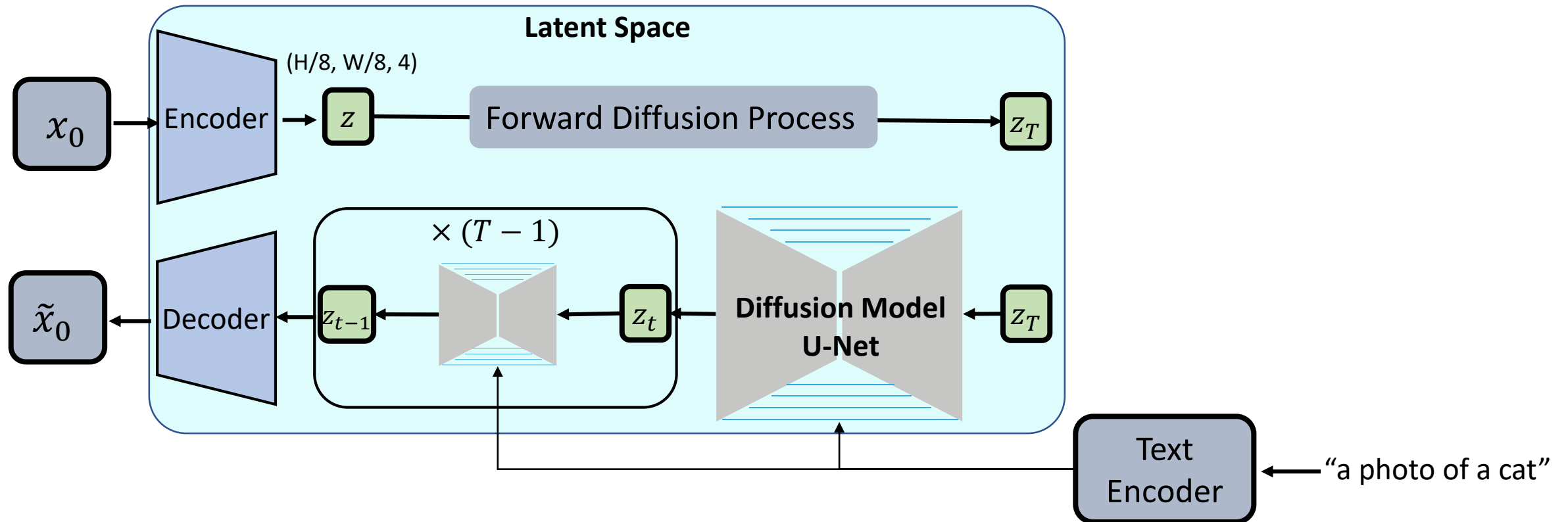
Fine



2048 x 2048

Does not fit in memory!

Latent Diffusion Models



Compress images with encoder \rightarrow diffusion steps \rightarrow decompress with a decoder

Everything as usual, except for training – it is done in **latent space**

Results



Summary

- A cross-modal network (CLIP) trained for discriminative tasks, can be used for generation as well.
- To generate images from text we need a generative prior:
 - VQ-GAN
 - StyleGAN
 - DIP
 - Diffusion Models

What we had today

1. CLIP-guided optimization:
 - VQ-GAN + CLIP
 - StyleCLIP
 - Text2LIVE
2. Diffusion Models + text
 - Text conditioning in Diffusion Models
 - Classifier (free) guidance
 - Latent Diffusion models

Questions?