



UiT The Arctic University of Norway

NORA summer school on multi-modal learning

Multi-modal Generation

Rwiddhi Chakraborty

UiT Machine Learning Group and Visual Intelligence

Schedule Today

- 10 – 11 : Generative AI - I
- 11 – 12 : Generative AI - II
- 12 – 13: Lunch
- 13 – 14: Multi-modal Generation
- 14 – 16: Group Project

In this talk

Conditioning

Control

Edits

Cool Stuff

In this talk

Conditioning

Control

Edits

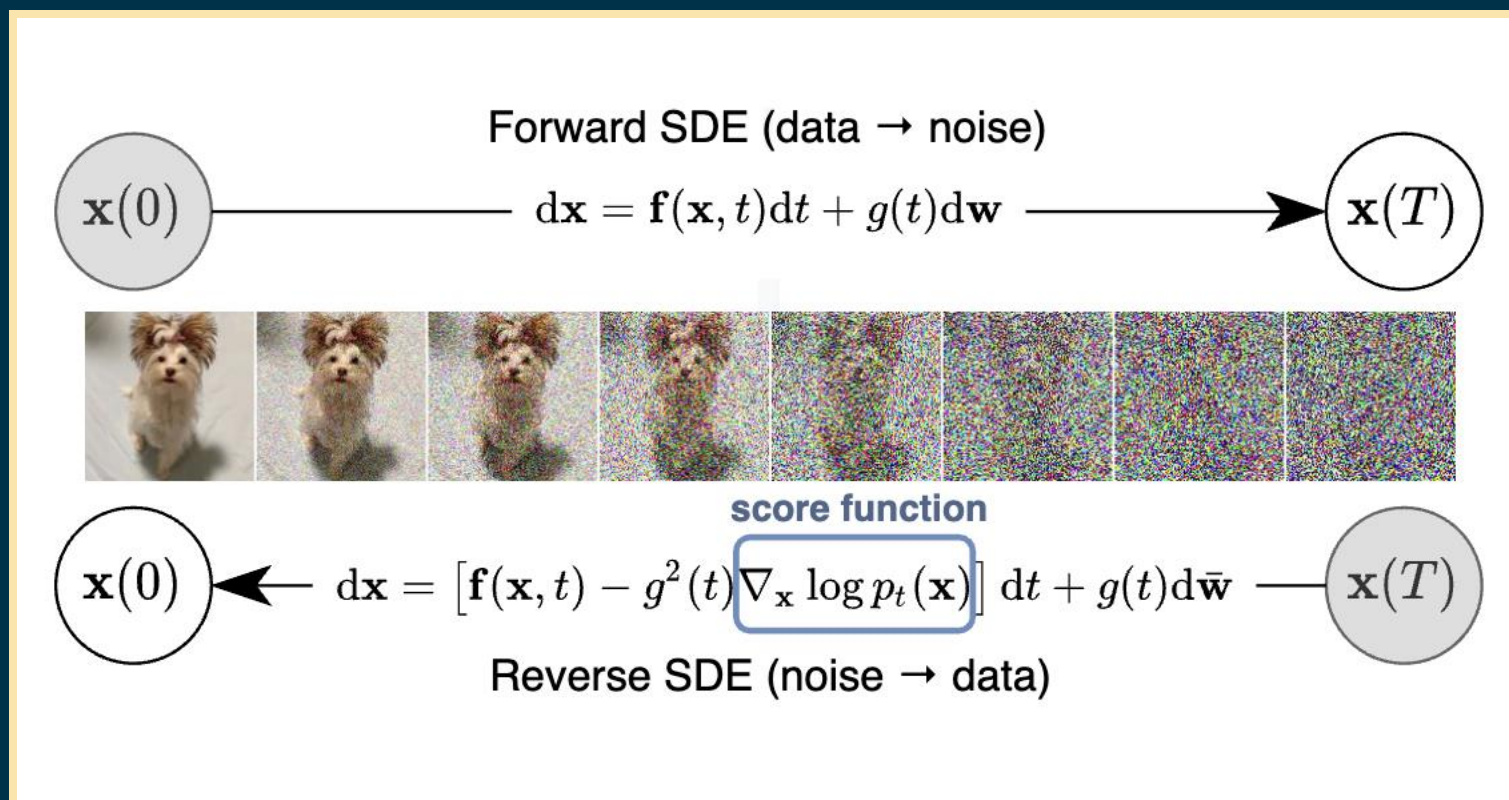
Cool Stuff

Generation: A noisy view

“Creating noise from data is easy, creating data from noise is generative modelling.”

Yang Song et al, Score-based generative modelling through stochastic differential equations.

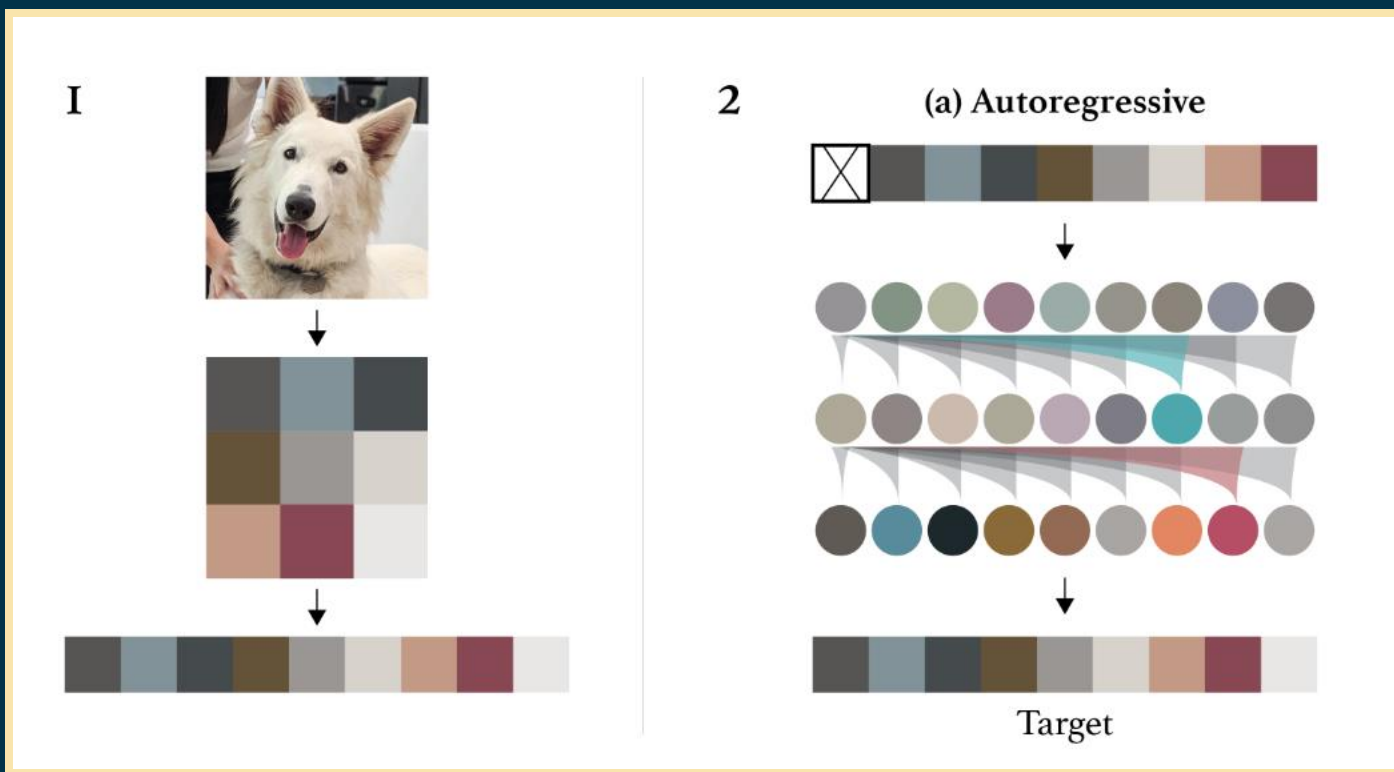
Generation: A noisy view



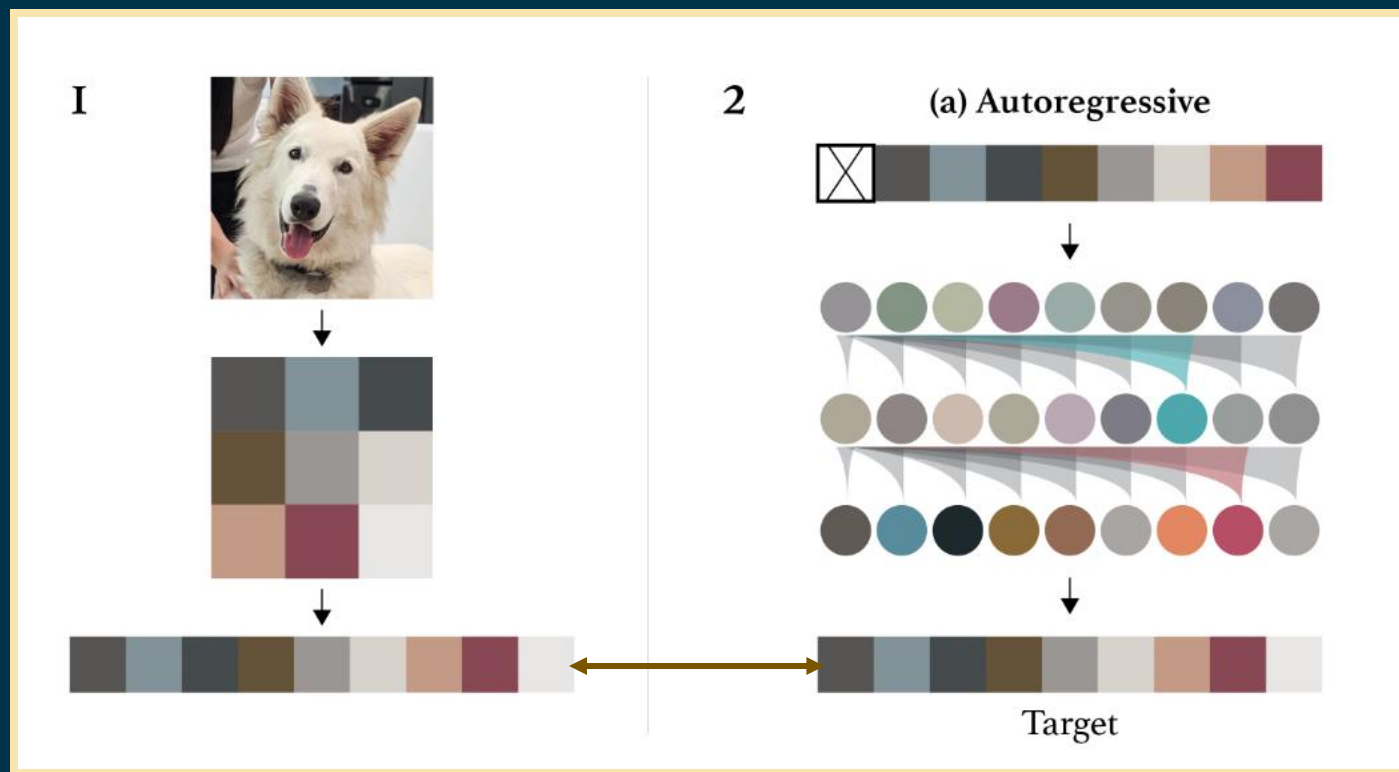
Generation: An autoregressive view

“....promising as our architecture uses a dense connectivity pattern which does not encode the 2D spatial structure of images yet is able to match and even outperform approaches which do....”

Generation: An autoregressive view



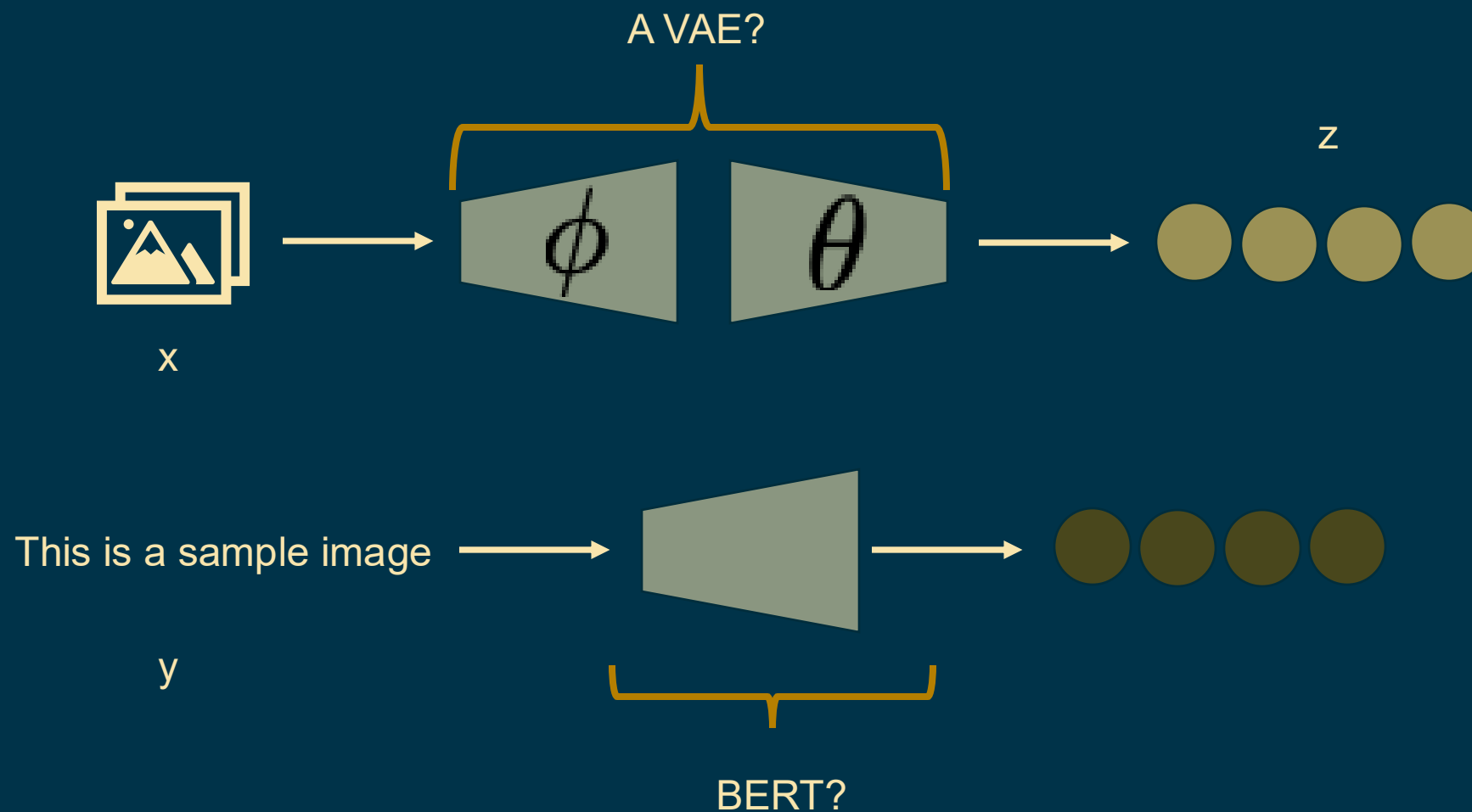
Generation: An autoregressive view



Generation

But how do we include another modality?

Conditioning: A first attempt



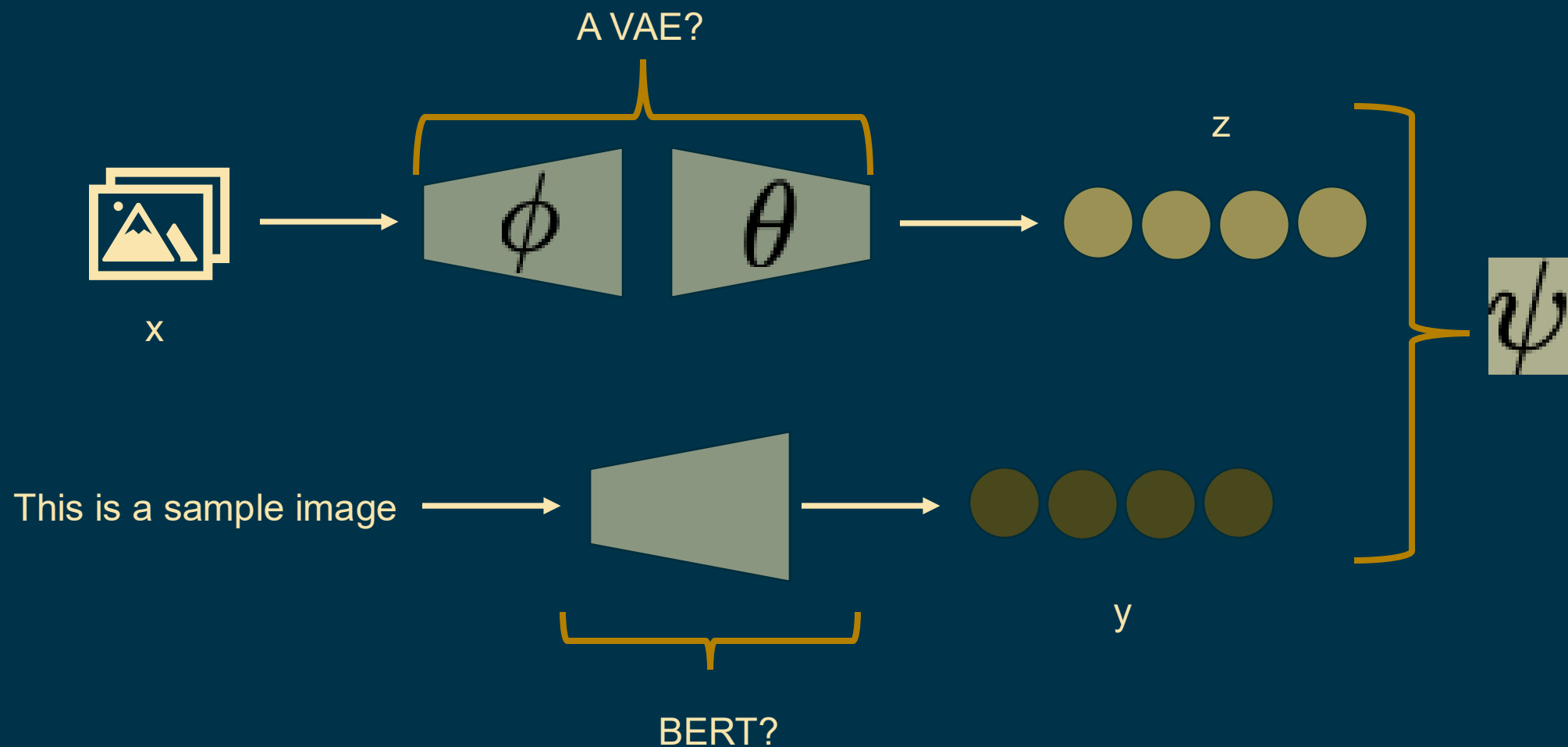
Conditioning: A first attempt

We have representations of the image and the text

Conditioning: A first attempt

We want to jointly model x , y , z

Conditioning: A first attempt



Conditioning: A first attempt

$$\ln p_{\theta, \psi}(x, y) \geq \mathbb{E}_{z \sim q_{\phi}(z|x)} (\underbrace{\ln p_{\theta}(x \mid y, z)}_{\text{VAE decoder}} - \beta D_{\text{KL}}(q_{\phi}(y, z \mid x) \parallel p_{\psi}(y, z)))$$

VAE decoder

Conditioning: A first attempt

$$\ln p_{\theta, \psi}(x, y) \geq \mathbb{E}_{z \sim q_{\phi}(z|x)} \left(\underbrace{\ln p_{\theta}(x \mid y, z)}_{\text{VAE decoder}} - \beta \underbrace{D_{\text{KL}}(q_{\phi}(y, z \mid x) \parallel p_{\psi}(y, z))}_{\text{VAE encoder}} \right)$$

Conditioning: A first attempt

$$\ln p_{\theta, \psi}(x, y) \geq \mathbb{E}_{z \sim q_{\phi}(z|x)} \left(\underbrace{\ln p_{\theta}(x \mid y, z)}_{\text{VAE decoder}} - \beta \underbrace{D_{\text{KL}}(q_{\phi}(y, z \mid x) \parallel p_{\psi}(y, z))}_{\text{VAE encoder Transformer}} \right)$$

Conditioning: A first attempt

$$\ln p_{\theta, \psi}(x, y) \geq \mathbb{E}_{z \sim q_{\phi}(z|x)} (\underbrace{\ln p_{\theta}(x | y, z)}_{\text{VAE decoder}} - \underbrace{\beta D_{\text{KL}}(q_{\phi}(y, z | x) || p_{\psi}(y, z))}_{\text{VAE encoder}})$$

Transformer

ELBO

The diagram illustrates the Evidence Lower Bound (ELBO) equation for a generative model. The equation is $\ln p_{\theta, \psi}(x, y) \geq \mathbb{E}_{z \sim q_{\phi}(z|x)} (\ln p_{\theta}(x | y, z) - \beta D_{\text{KL}}(q_{\phi}(y, z | x) || p_{\psi}(y, z)))$. Brackets below the equation identify its parts:

- $\ln p_{\theta}(x | y, z)$ is labeled "VAE decoder".
- $\beta D_{\text{KL}}(q_{\phi}(y, z | x) || p_{\psi}(y, z))$ is labeled "VAE encoder".
- $p_{\psi}(y, z)$ is labeled "Transformer".
- A large bracket under the entire right-hand side of the inequality is labeled "ELBO".

Conditioning: A first attempt

Two-stage training

First, train the VAE

Initial prior a uniform distribution over a K-codebook

=> Maximise ELBO wrt θ and ϕ

Conditioning: A first attempt

Two-stage training

First, train the VAE

Initial prior a uniform distribution over a K-codebook

=> Maximise ELBO wrt θ and ϕ

Next, train the transformer

=> Maximise ELBO wrt ψ

Text and image concatenated as a single stream of data

Conditioning: A first attempt



Conditioning: A second attempt

Are there some issues with the autoregressive approach?

Conditioning: A second attempt

“....promising as our architecture uses a dense connectivity pattern which does not encode the 2D spatial structure of images yet is able to match and even outperform approaches which do....”

Conditioning: A second attempt

“....promising as our architecture uses a dense connectivity pattern which does not encode the 2D spatial structure of images yet is able to match and even outperform approaches which do”

Conditioning: A second attempt

We need inductive bias

Conditioning: A second attempt

Are there some issues with the autoregressive approach?

Conditioning: A second attempt

High-dimensional modelling in pixel-space is extremely inefficient

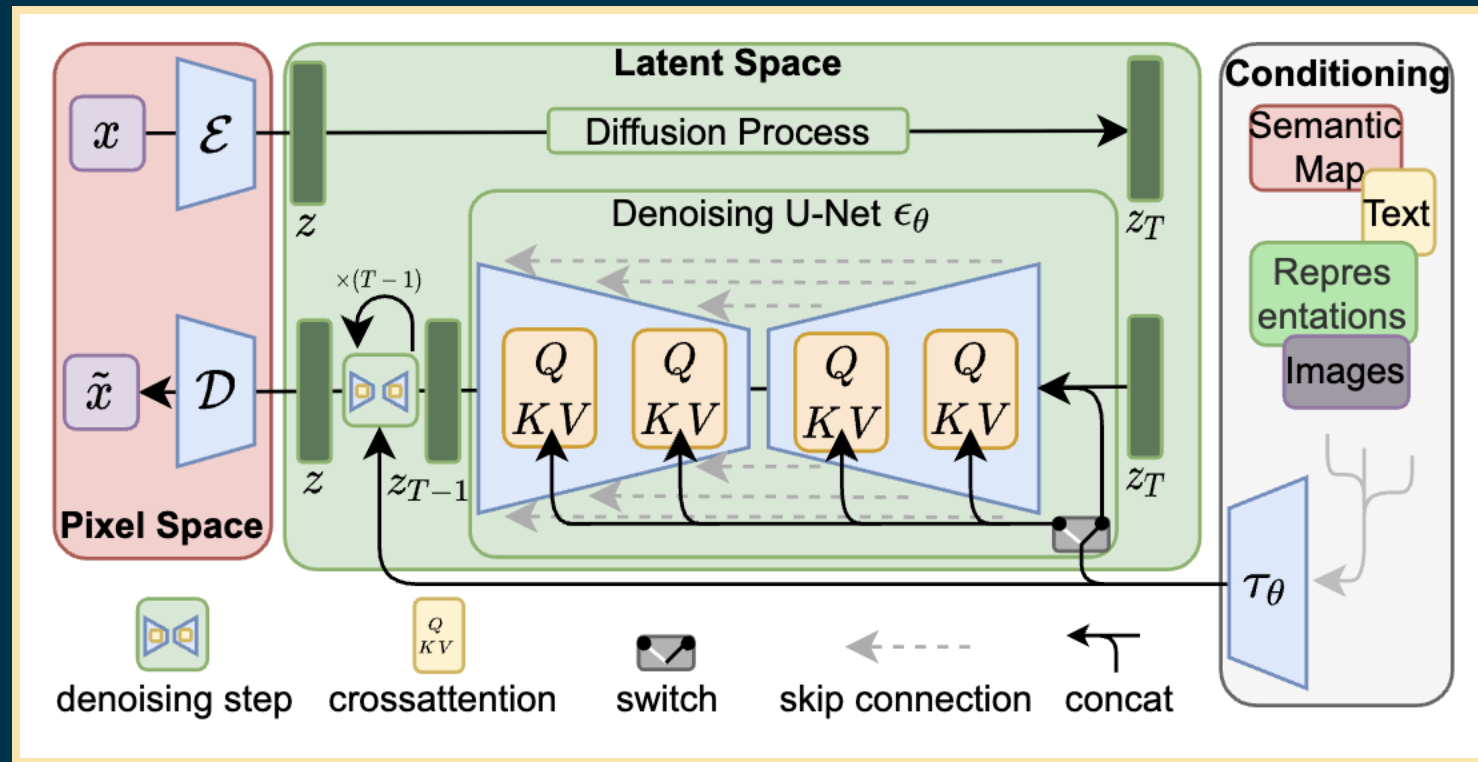
Conditioning: A second attempt

U-Net and 2D convolutions to preserve spatial structure

Compressed latent space with a VAE backbone

Cross-attention between text and image embeddings

Conditioning: A second attempt



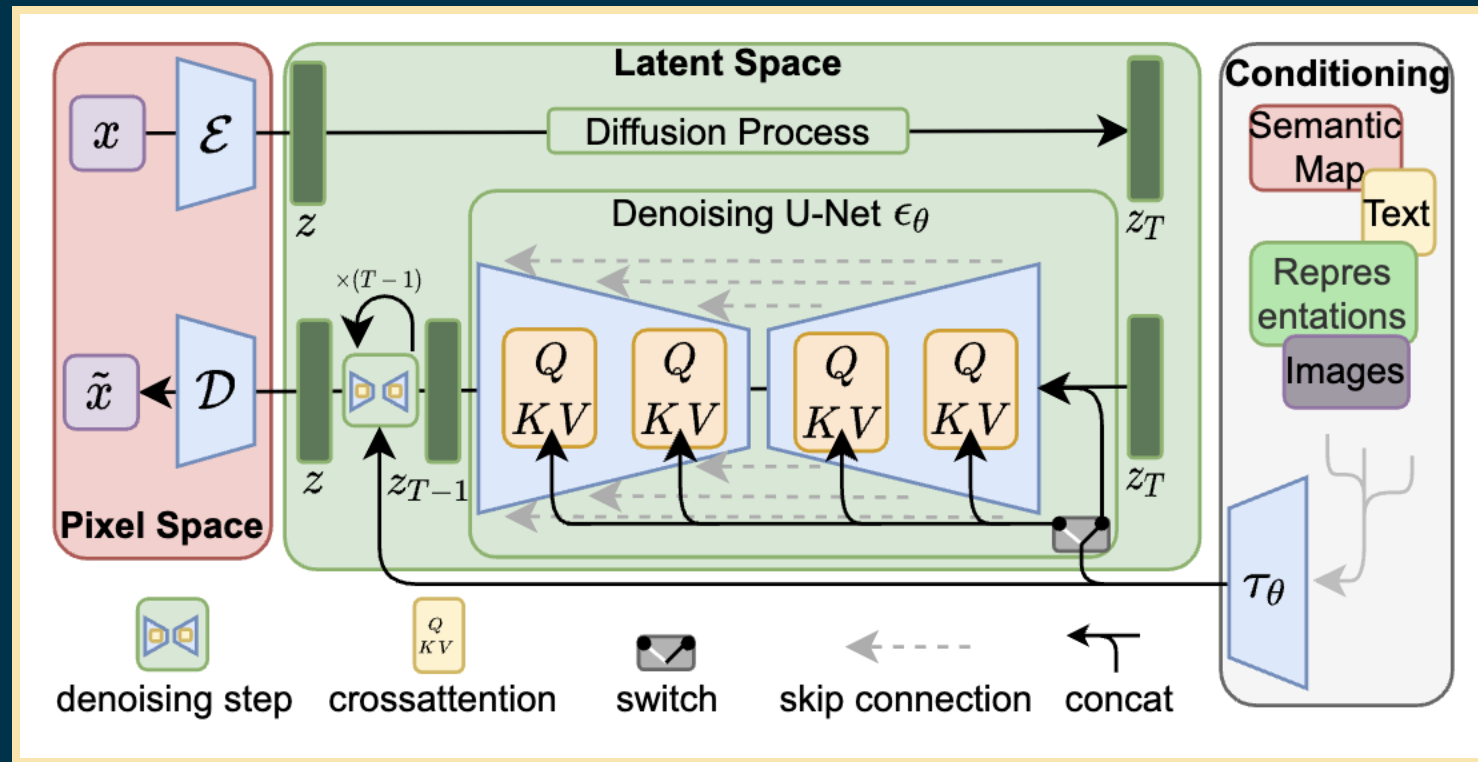
Conditioning: A second attempt

DDIM
↓

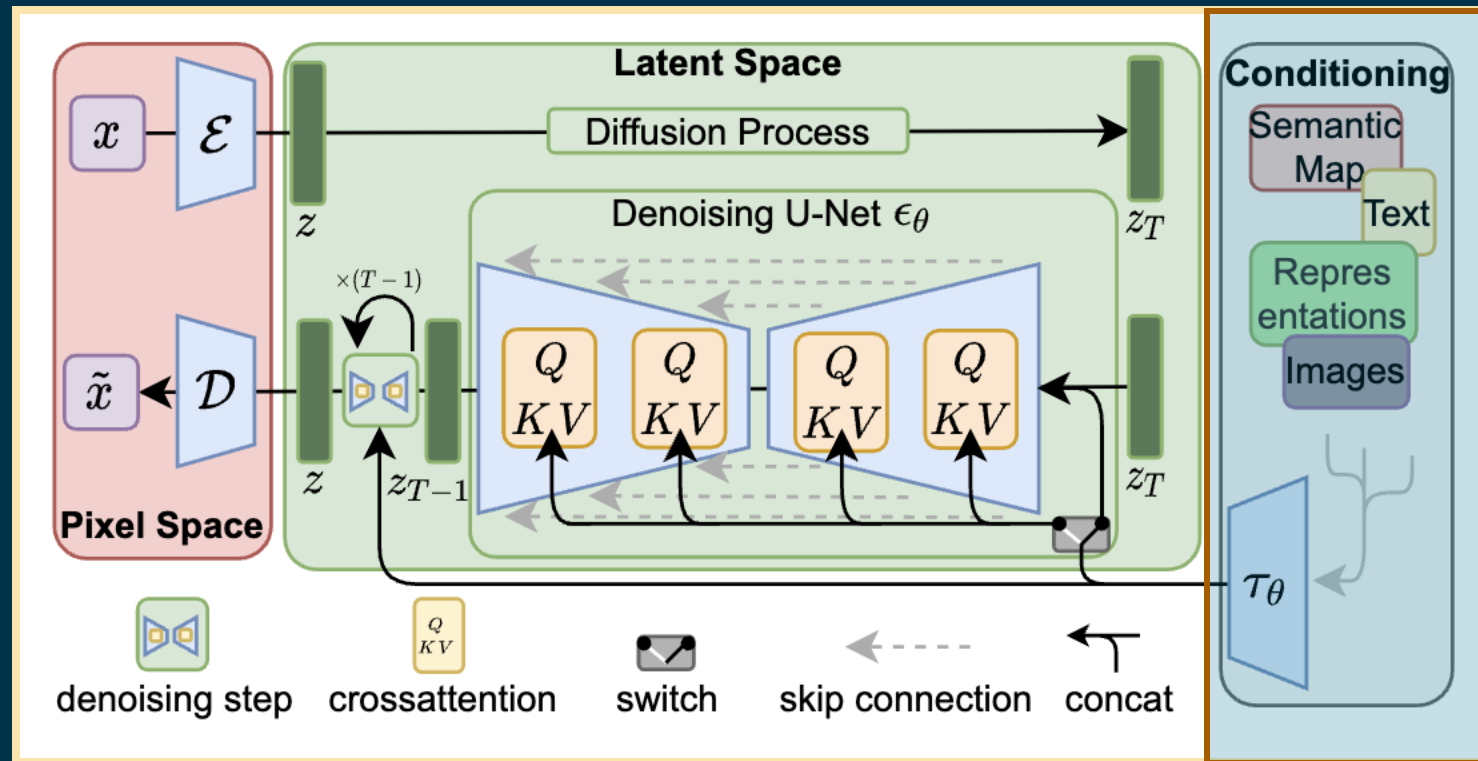
$$\mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|_2^2 \right]$$

↑ ↑
Denoising U-Net Sparse transformer

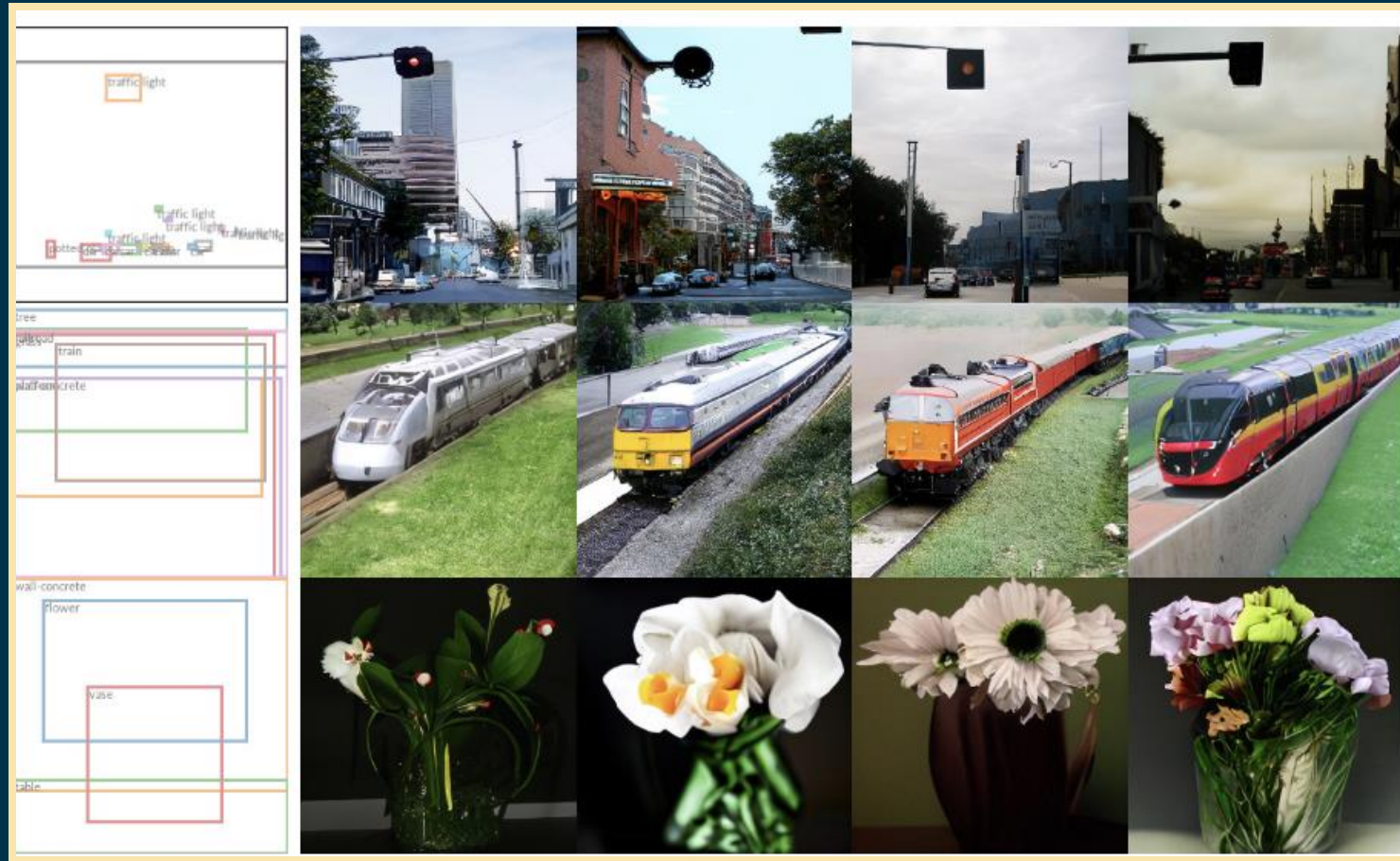
Conditioning Zoo



Conditioning Zoo



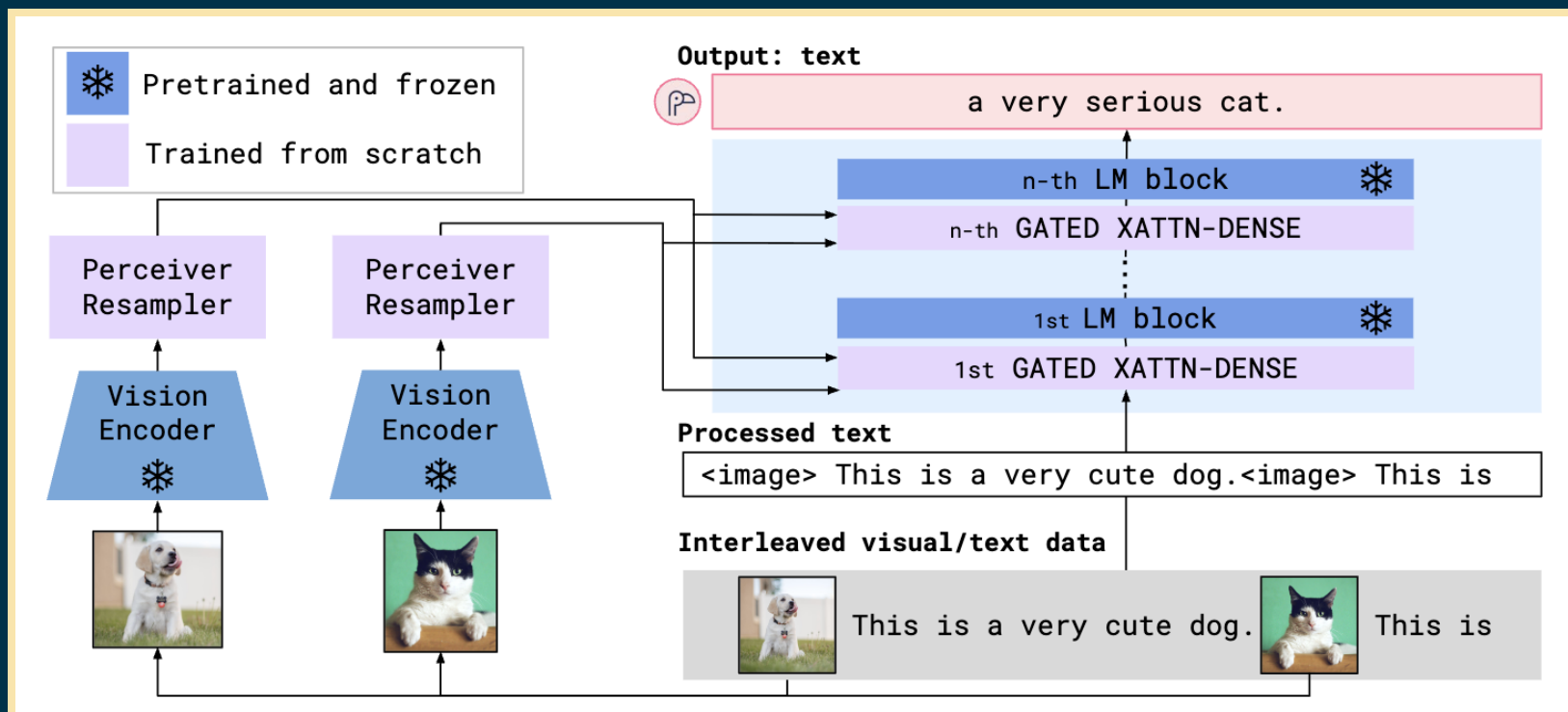
Beyond Text



Beyond Text



Generative vs Discriminative: A blurry line



In this talk

Conditioning

Control

Edits

Cool Stuff

Control: A form of fine-tuning

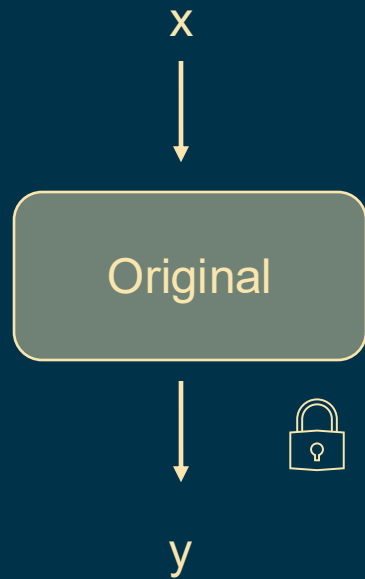
Given a pretrained text-to-image model

Can we generate new, spatially localized, task specific images?

Control: A form of fine-tuning



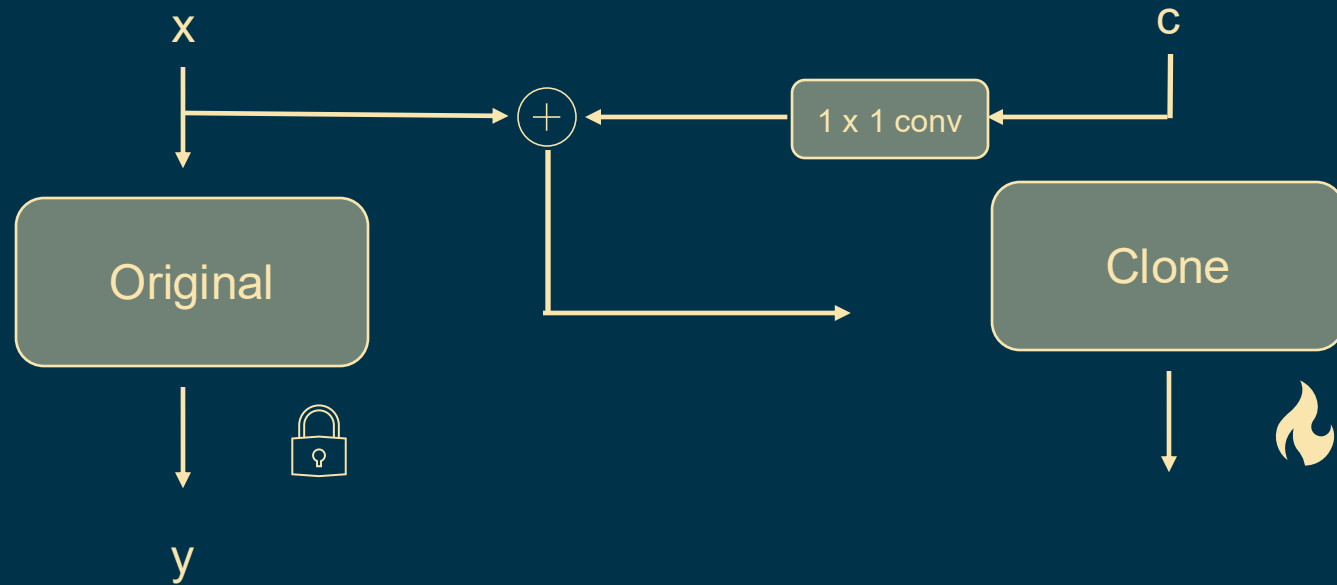
Control: A residual on a clone



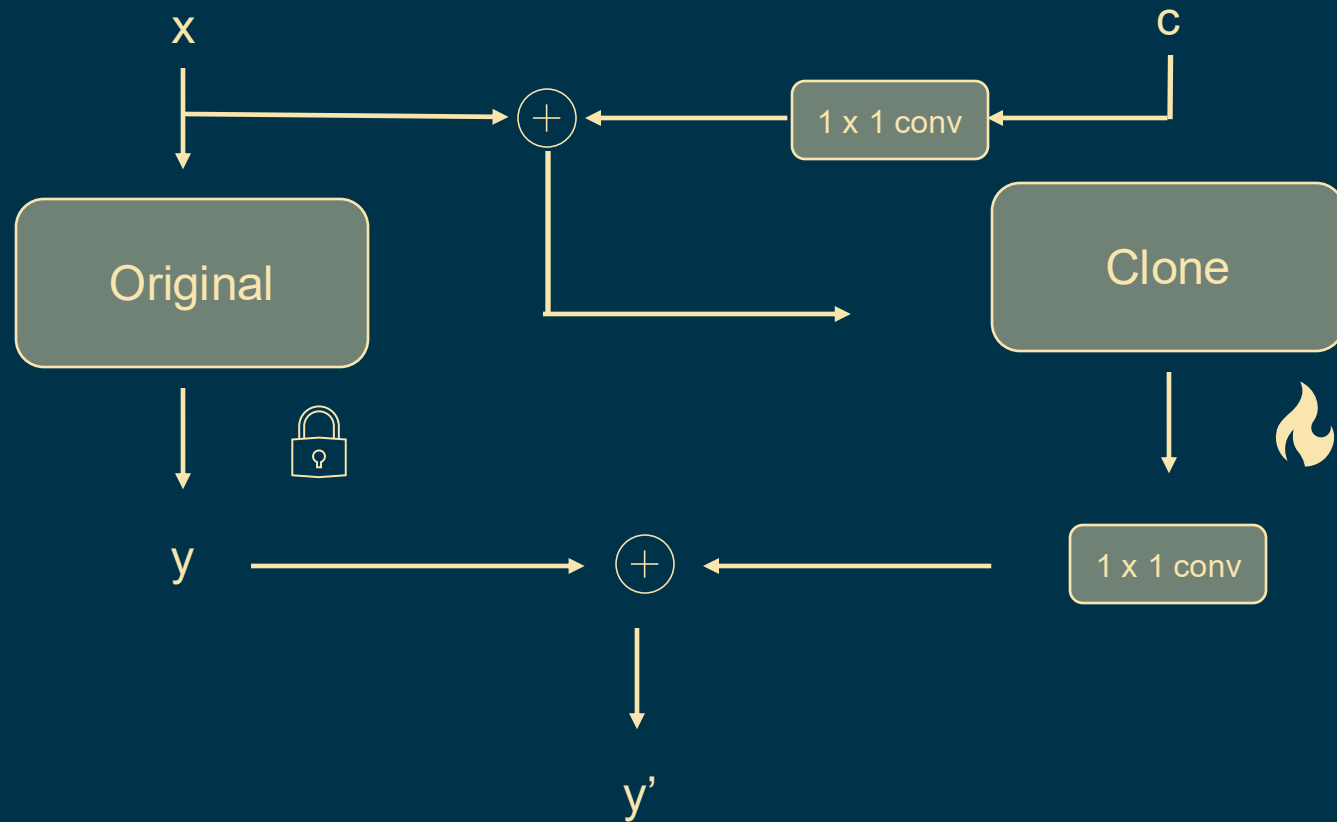
Control: A residual on a clone



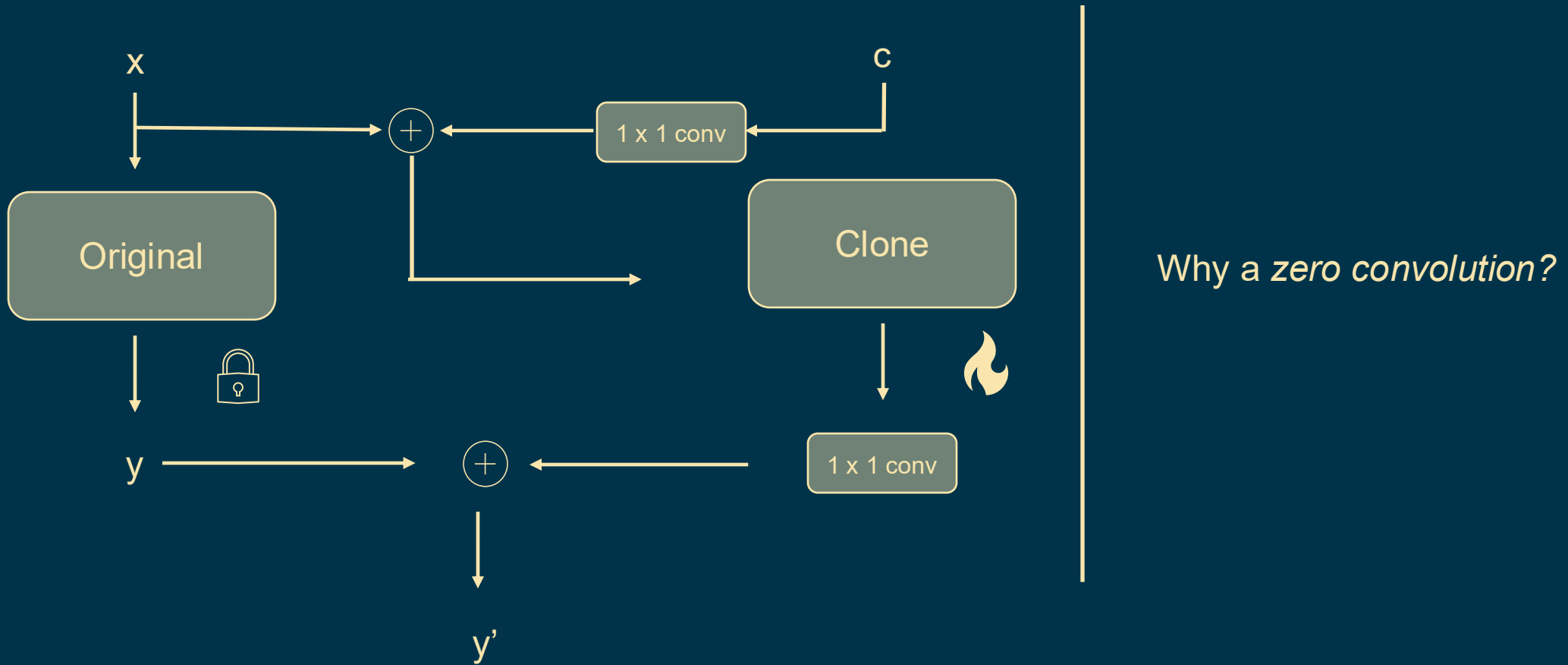
Control: A residual on a clone



Control: A residual on a clone



Control: A residual on a clone

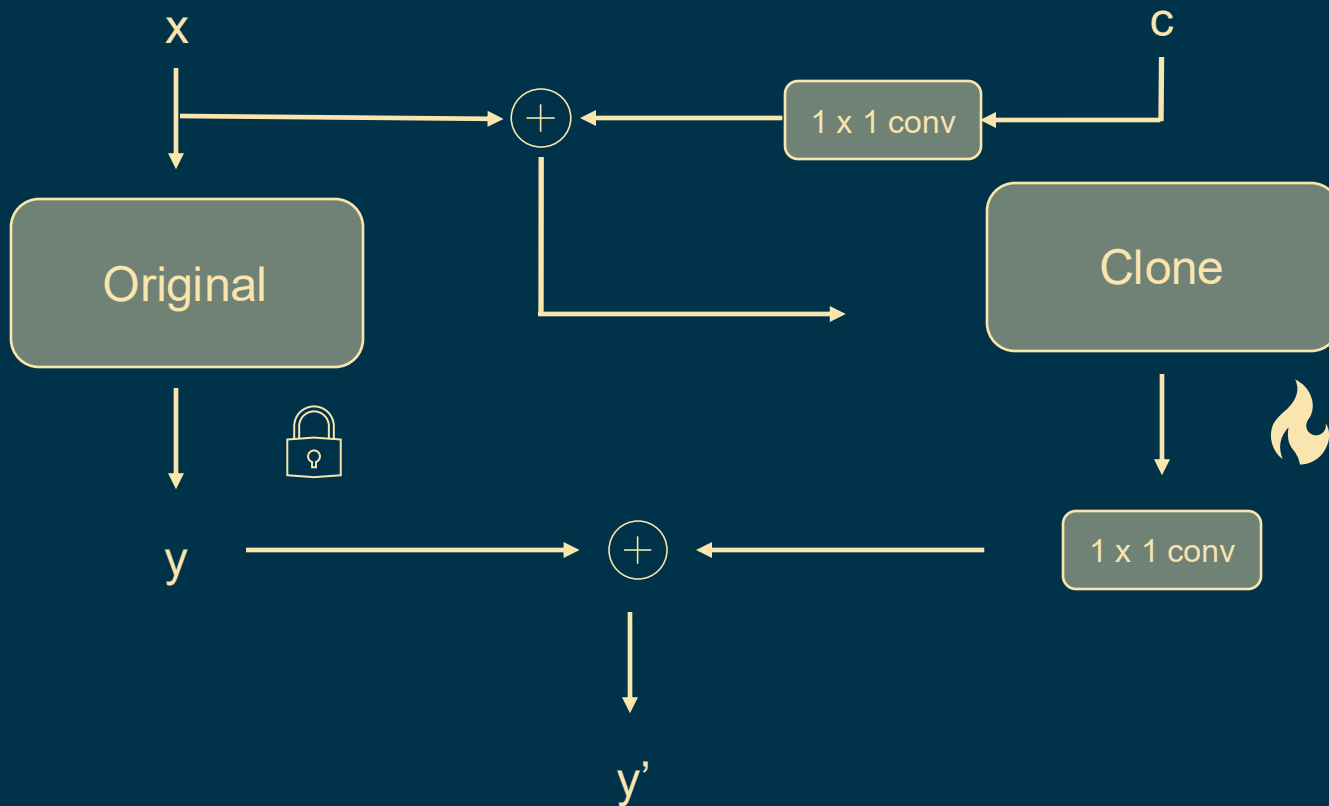


The Zero Convolution

$$y_c = \mathcal{F}(x; \Theta) + \mathcal{Z}(\mathcal{F}(x + \mathcal{Z}(c; \Theta_{z1}); \Theta_c); \Theta_{z2})$$

What happens when zero?

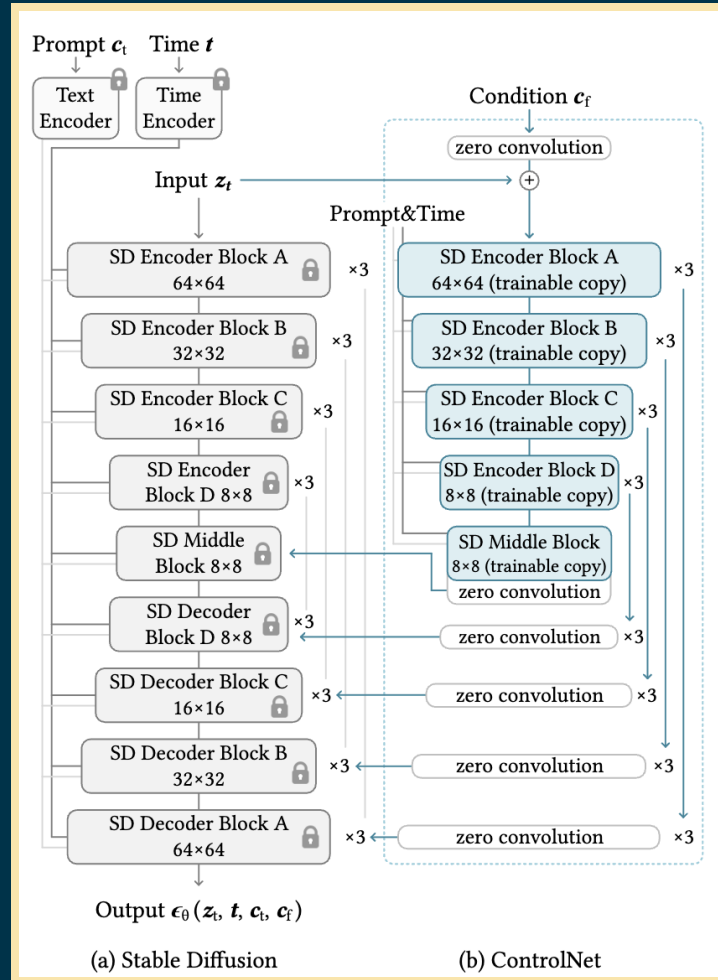
Control: A residual on a clone



Why a *zero convolution*?

In the first forward pass, keep things unchanged

ControlNet



ControlNet

$$\mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|_2^2 \right]$$

Before

ControlNet

$$\mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|_2^2 \right]$$

Before

$$\mathbb{E}_{z_0, t, c_t, c_f, \epsilon \sim \mathcal{N}(0,1)} \left[\|\epsilon - \epsilon_{\theta}(z_t, \mathbf{t}, \mathbf{c}_t, c_f)\|_2^2 \right]$$

After

ControlNet

$$\mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|_2^2 \right]$$

Before

$$\mathbb{E}_{z_0, t, c_t, c_f, \epsilon \sim \mathcal{N}(0,1)} \left[\|\epsilon - \epsilon_{\theta}(z_t, \mathbf{t}, \mathbf{c}_t, c_f)\|_2^2 \right]$$

After

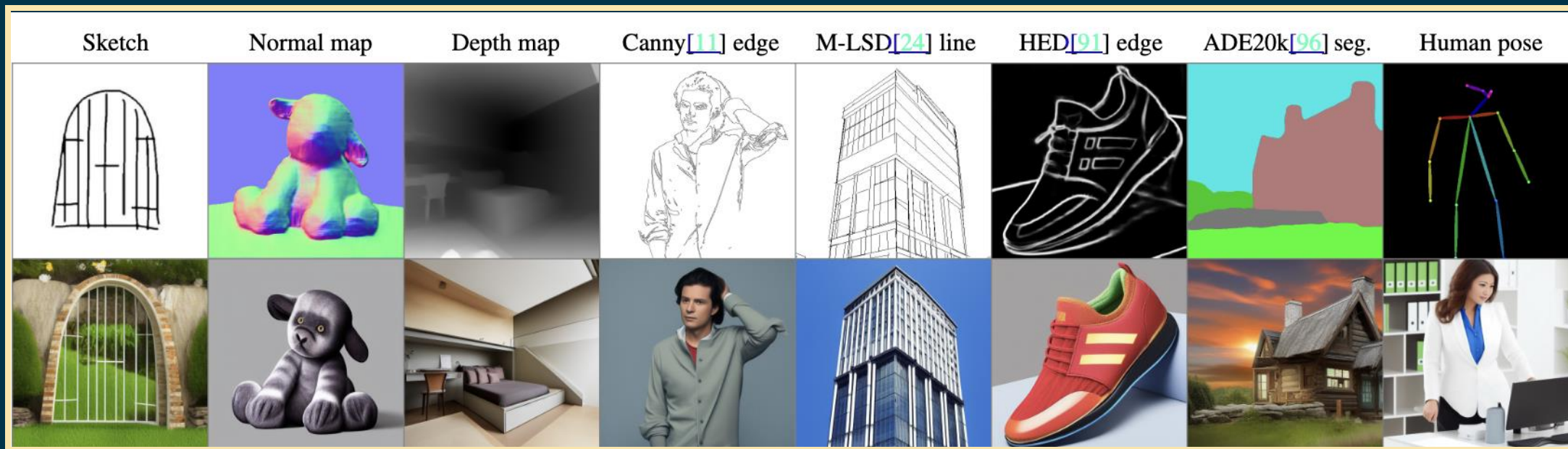
Text prompt



Task vector



ControlNet: Text not a necessity!



$$\mathbb{E}_{z_0, t, c_t, c_f, \epsilon \sim \mathcal{N}(0,1)} \left[\|\epsilon - \epsilon_\theta(z_t, t, \mathbf{c}_t, c_f)\|_2^2 \right]$$

Zero!

In this talk

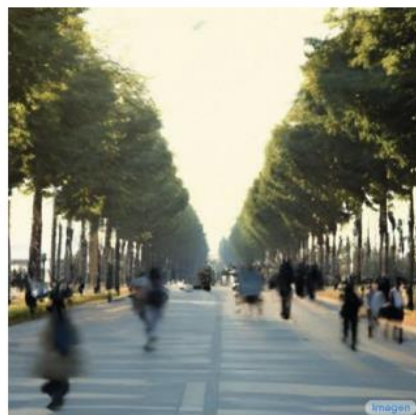
Conditioning

Control

Edits

Cool Stuff

Text driven image editing



“The boulevards are crowded today.”

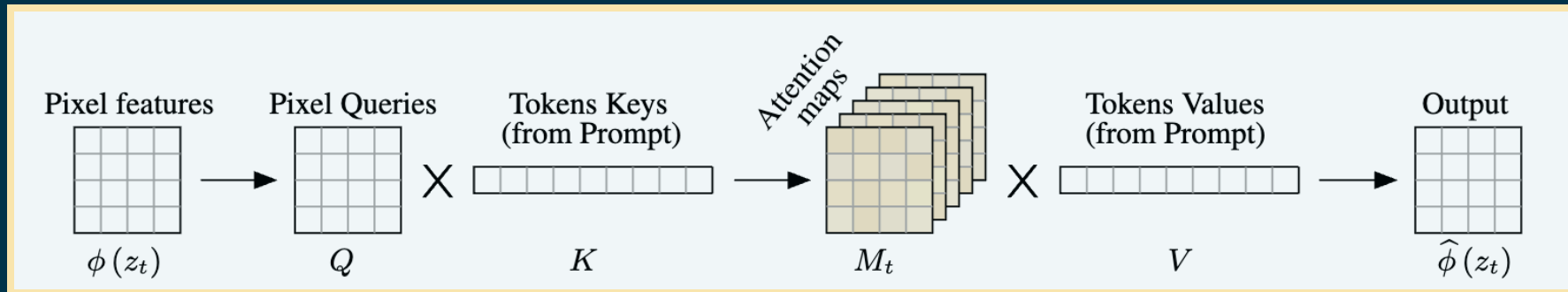


“Photo of a cat riding on a bicycle.”

~~bicycle~~
car

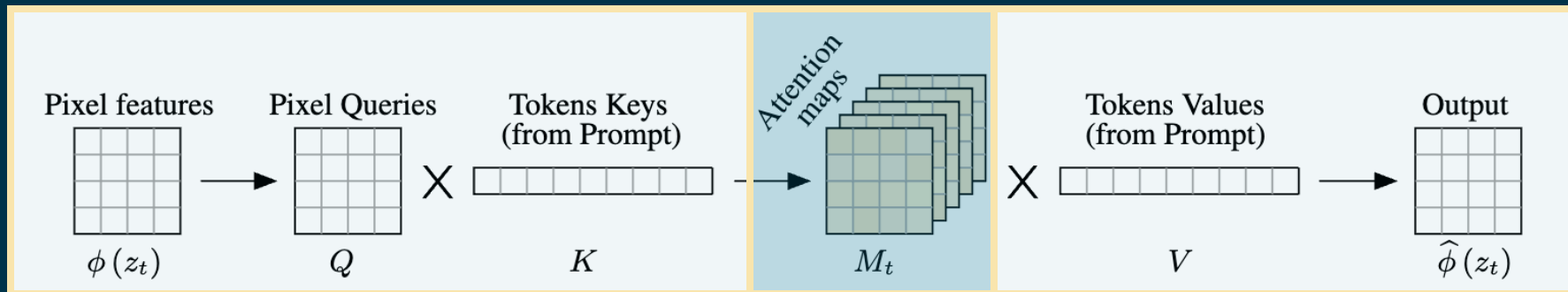
Text driven image editing

The cross-attention layers hold the information we need



Text driven image editing

The cross-attention layers hold the information we need



Text driven image editing

The cross-attention layers hold the information we need

A photo of a cat on a **bicycle** 

Prompt P

Attention map M

Text driven image editing

The cross-attention layers hold the information we need

A photo of a cat on a **bicycle**



A photo of a cat on a **car**

Prompt P

Prompt P*

Attention map M

Attention map M*

Text driven image editing

The cross-attention layers hold the information we need

Preserve original layout and geometry

A photo of a cat on a **bicycle**



A photo of a cat on a **car**

Prompt P

Prompt P*

Attention map M

Attention map M*

Text driven image editing

The cross-attention layers hold the information we need

A photo of a cat on a **bicycle**

Prompt P

Attention map M



A photo of a cat on a **car**

Prompt P*

Attention map M*

Text driven image editing

The cross-attention layers hold the information we need

A photo of a cat on a **bicycle**

Prompt P **✗**

Attention map M **✓**



A photo of a cat on a **car**

Prompt P* **✓**

Attention map M* **✗**

Text driven image editing

The cross-attention layers hold the information we need

A photo of a cat on a **bicycle**

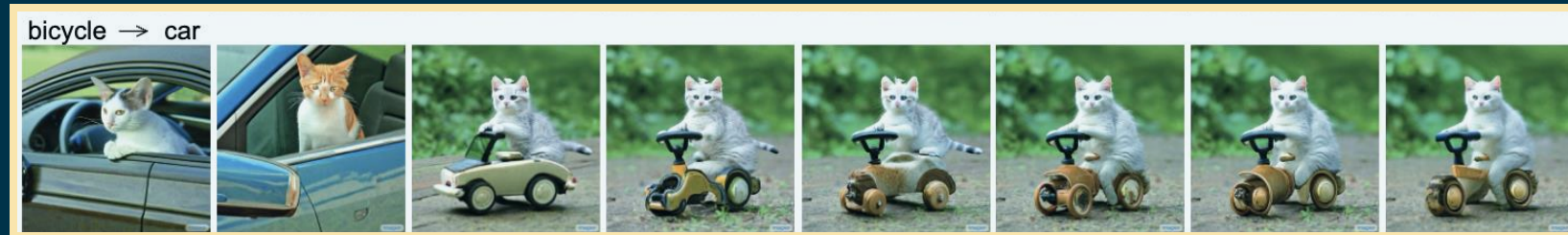
Prompt P ✗

Attention map M ✓

A photo of a cat on a **car**

Prompt P* ✓

Attention map M* ✗



Text driven image editing – Phrases

The cross-attention layers hold the information we need

A photo of a car on the side of the street



A photo of a **sports car** on the side of the street

Prompt P

Attention map M

Prompt P*

Attention map M*

Text driven image editing – Phrases

The cross-attention layers hold the information we need

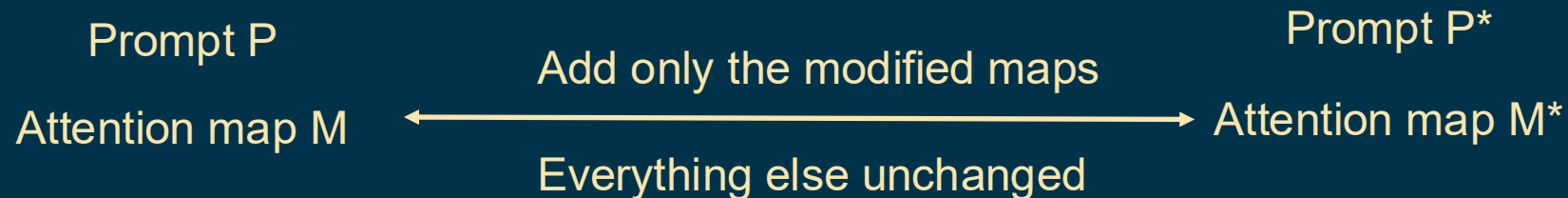
A photo of a car on the side of the street \longrightarrow A photo of a **sports car** on the side of the street



Text driven image editing – Phrases

The cross-attention layers hold the information we need

A photo of a car on the side of the street \longrightarrow A photo of a **sports car** on the side of the street



Text driven image editing

The cross-attention layers hold the information we need



Text driven image editing

How about instructions?

Text driven image editing

Instead of:

A photo of the Eiffel Tower  A photo of the Eiffel tower with fireworks in the sky

Text driven image editing

Instead of:

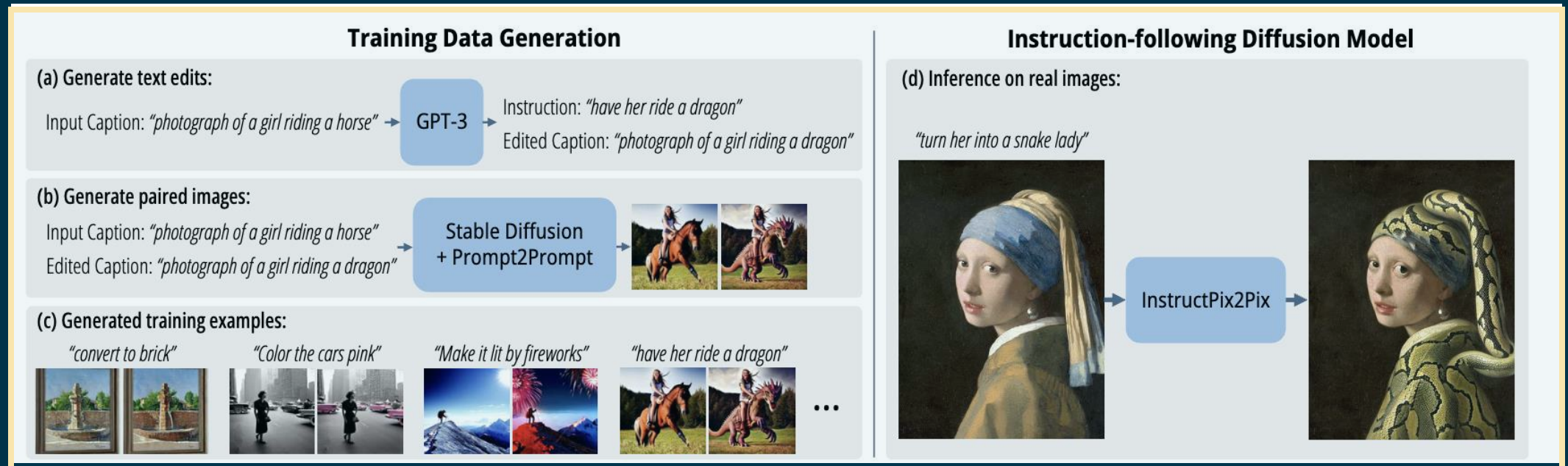
A photo of the Eiffel Tower  A photo of the Eiffel tower with fireworks in the sky

How about:

Add fireworks to the sky

Text driven image editing

Instruction following as supervised learning



In this talk

Conditioning

Control

Edits

Cool Stuff

Cool Stuff

Think about all of this as representation learning

We can use representations to do *anything*

Cool Stuff

Think about all of this as representation learning

So why not couple generative and discriminative models?

Cool Stuff

Recall that CLIP can classify images into labels

We now know that T2I models learn rich representations of world concepts

Cool Stuff

Recall that CLIP can classify images into labels



We now know that T2I models learn rich representations of world concepts

Open World Segmentation

Input Image



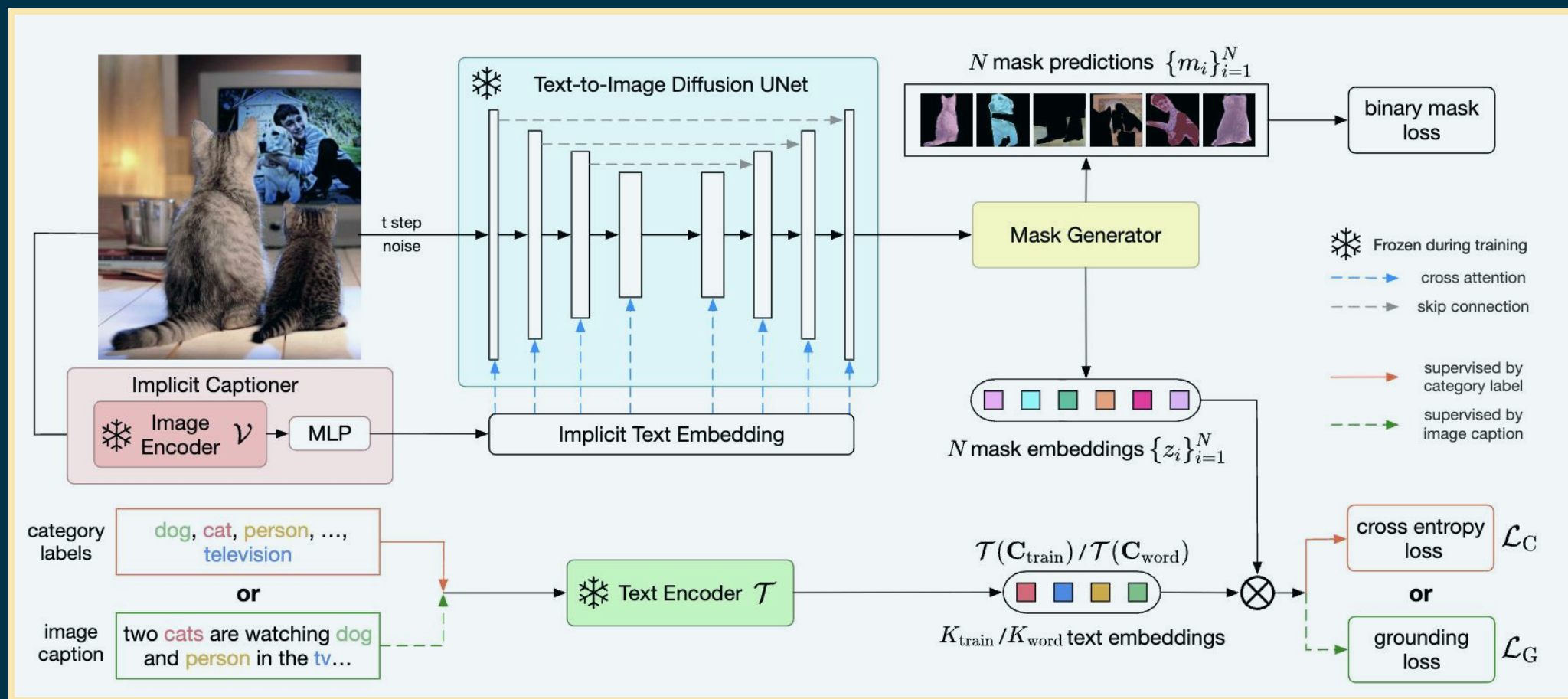
K-Means Clustering of Frozen Diffusion Features



Open-Vocabulary Panoptic Segmentation Prediction from ODISE



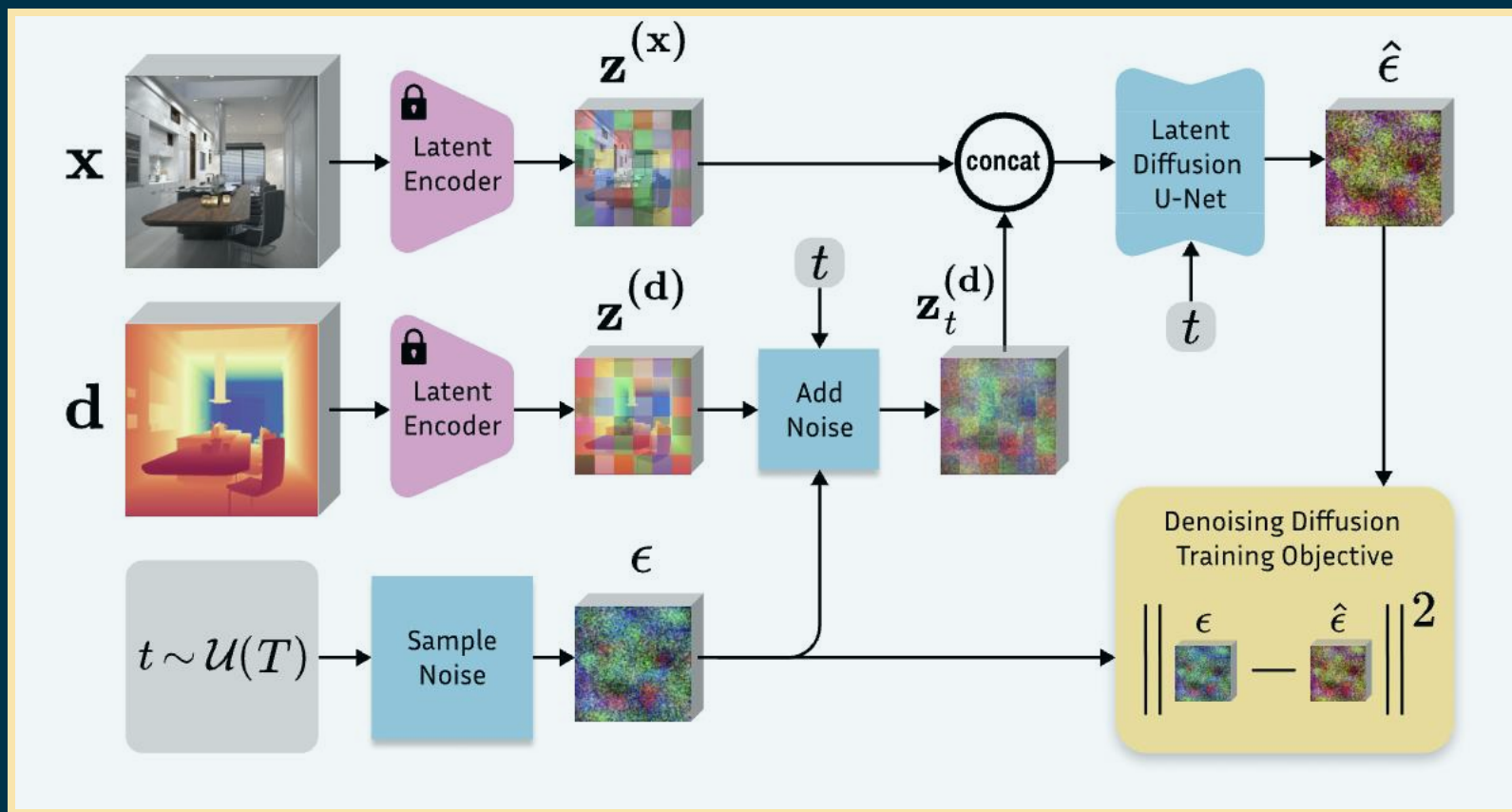
Open World Segmentation



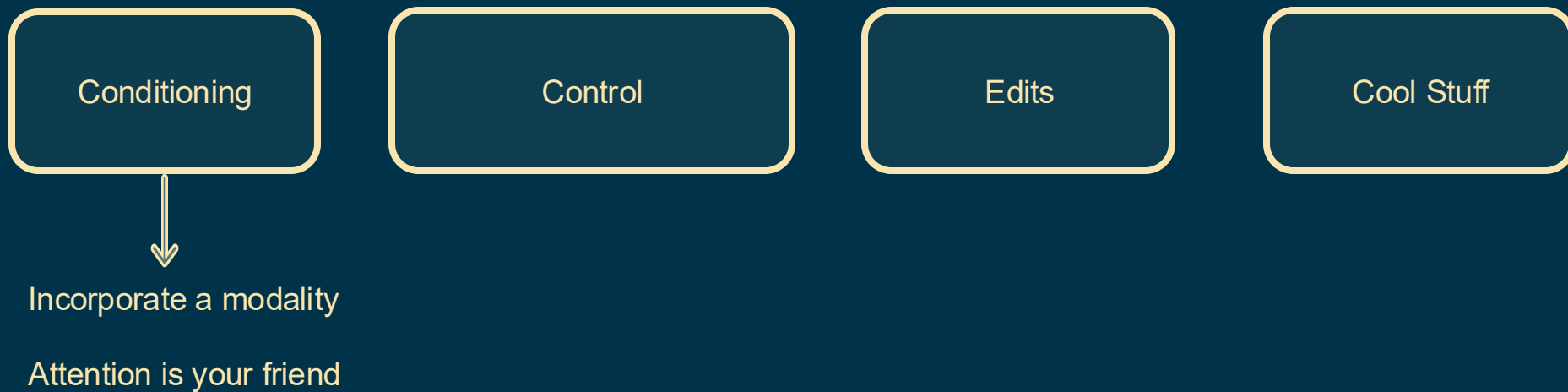
Depth Estimation



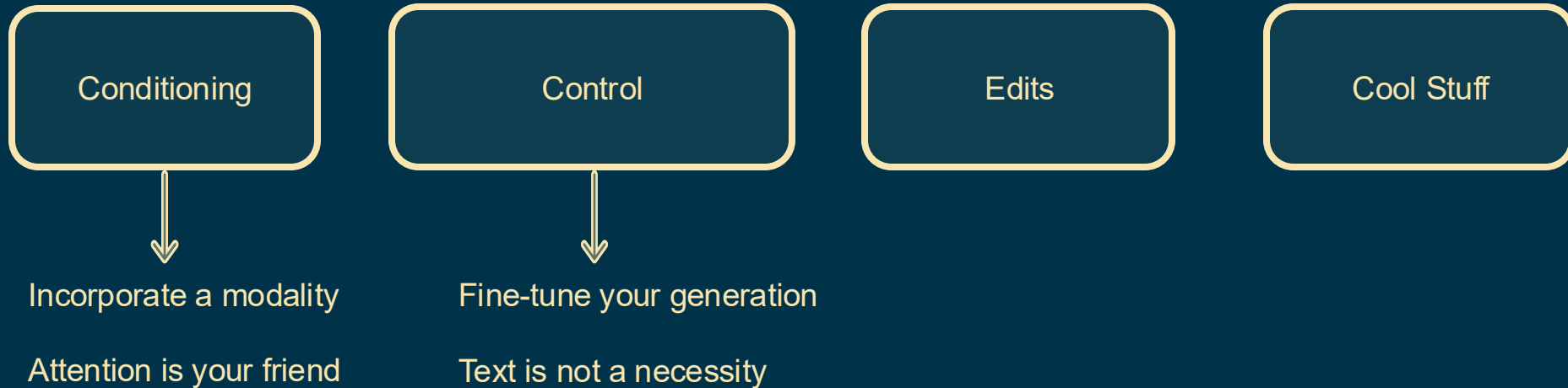
Depth Estimation



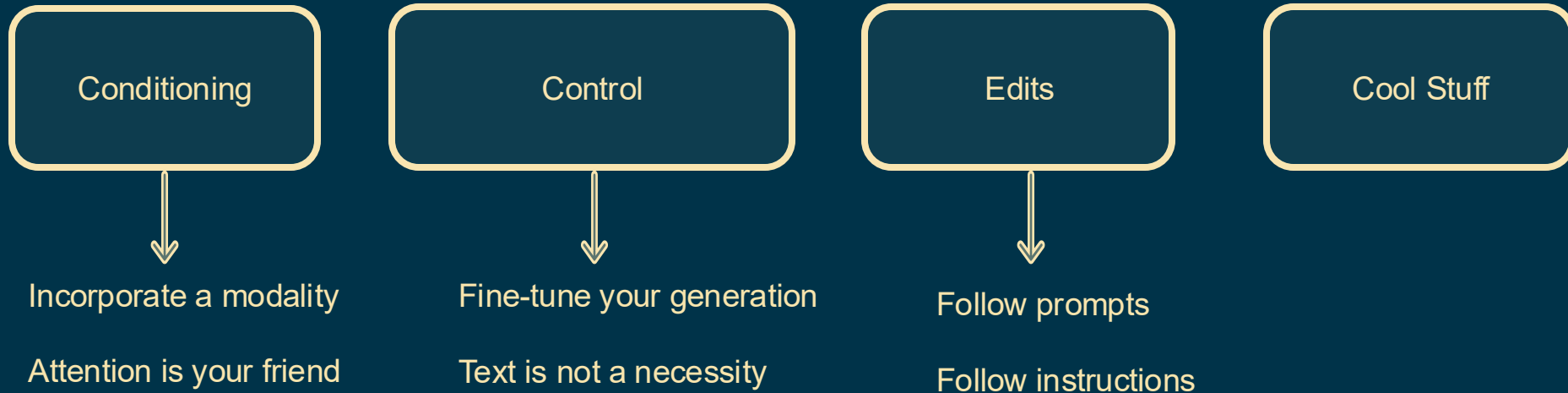
Summary



Summary



Summary



Summary

