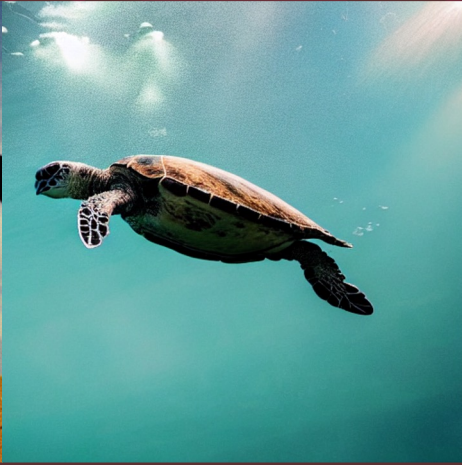# Stable Diffusion

Dongliang Guo, Jacobi Coleman
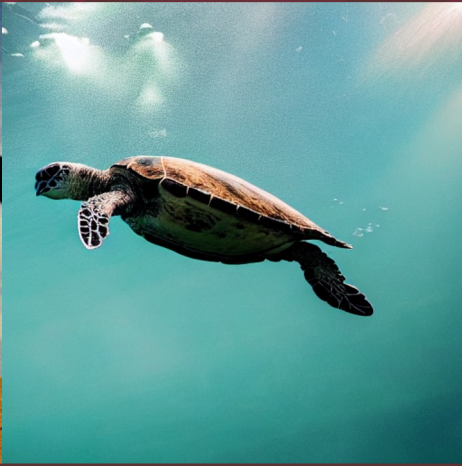
# Outline

- Introduction

- Motivation (related work)

- Problem Definition

- Methodology

- Results

- Future Direction

# What's the deal with all these pictures?

These pictures were generated by **Stable Diffusion,**
a recent diffusion generative model.

Along with other things, It can turn text prompts (e.g. "an astronaut riding a horse") into images.

# What makes this so important?

Allows for more creativity to be expressed without the confines of human physical capabilities.



"Multiple synapses firing around the brain"

# Why should we care?



"a lovely cat running in the desert in Van Gogh style, trending art."

Could be a model of imagination.

Similar techniques could be used to generate
any number of things (e.g. neural data).

It's *cool!*

"Batman eating pizza in a diner"

# How does it work?

It's complicated…
but here's the high-level idea.
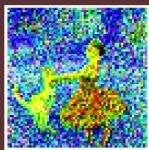
# What do we need?

Example pictures of people




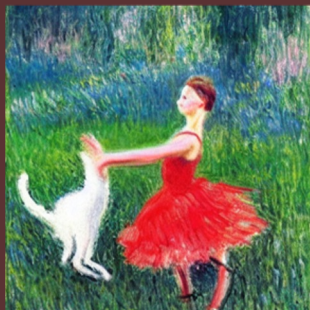


1. Method of learning to generate new stuff given many examples

# What do we need?

## 2. Way to link text and images

"cool professor person"

## 3. Way to compress images
(for speed in training and generation)

z[0:3, :, :]

# What do we need?

4. Way to add in good image-related inductive biases…

…since when you're generating something new, you need a way to safely go beyond the images you've seen before.

# What do we need?

1. Method of learning to generate new stuff

**Forward/reverse diffusion**

2. Way to link text and images

**Text-image representation model**

3. Way to compress images

**Autoencoder**

4. Way to add in good inductive biases

**U-net Architecture + 'attention'**

Making a 'good' generative model is about making all these parts work together well!

# Stable Diffusion in Action

"A mecha robot in a favela in expressionist style"

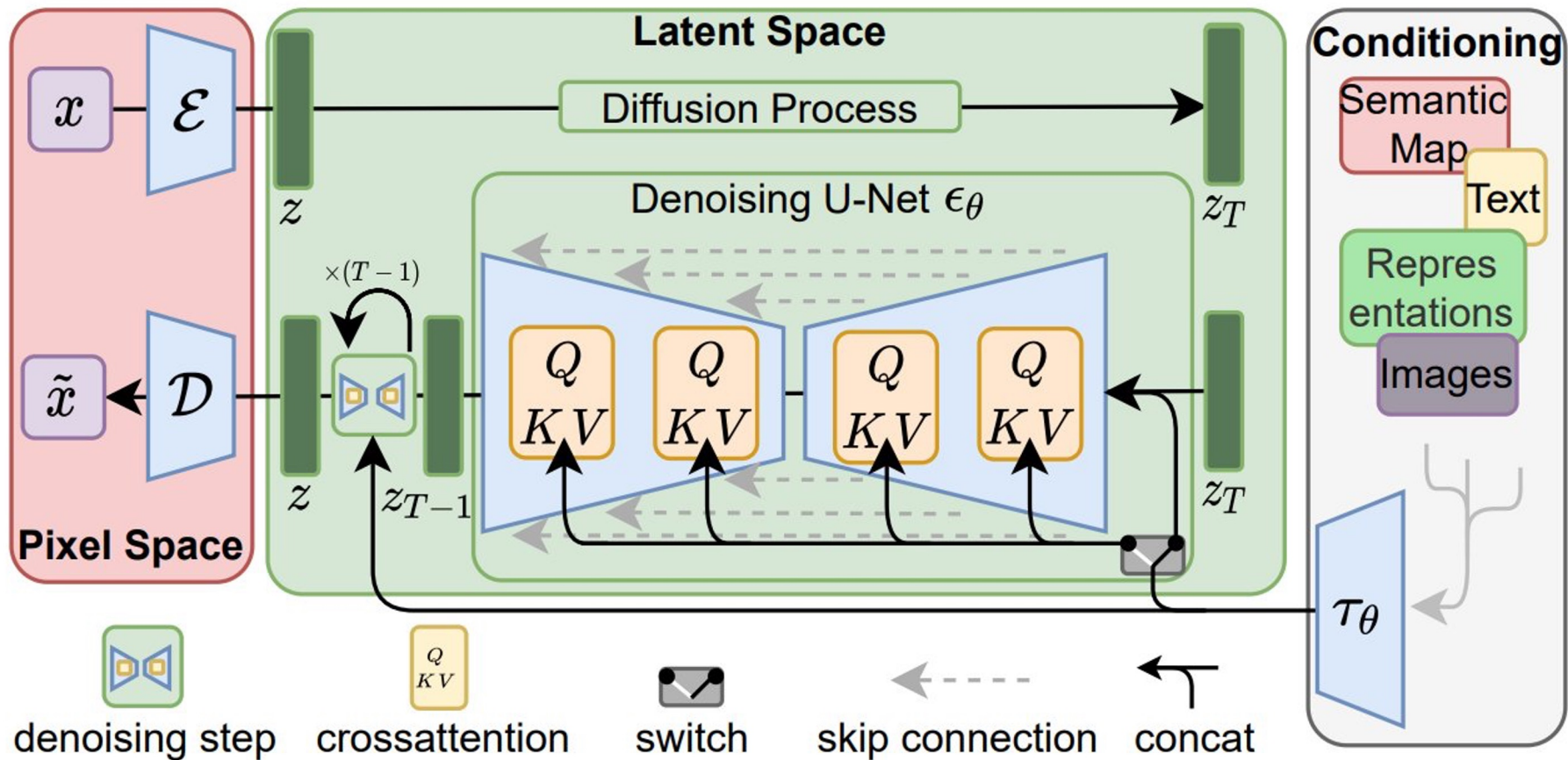# Cartoon with StableDiffusion + Cartoon

# Some Resources

- Diffusion model in general
    - What are Diffusion Models? | Lil'Log
    - Generative Modeling by Estimating Gradients of the Data Distribution | Yang Song

- Stable diffusion
    - Annotated & simplified code: U-Net for Stable Diffusion (labml.ai)
    - Illustrations: The Illustrated Stable Diffusion – Jay Alammar

- Attention & Transformers
    - The Illustrated Transformer

# What is the problem?

Training such a model requires massive computational resources only available to a small fraction of the field, and leaves a huge carbon footprint.

Secondly, evaluating an already trained model is also expensive in time and memory, since the same model architecture must run sequentially for a large number of steps .
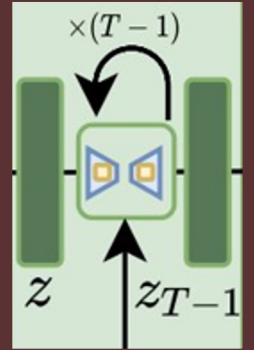
A goal of this research is to lower the computational demands of training diffusion models towards high-resolution image synthesis.

# Principle of Diffusion Models
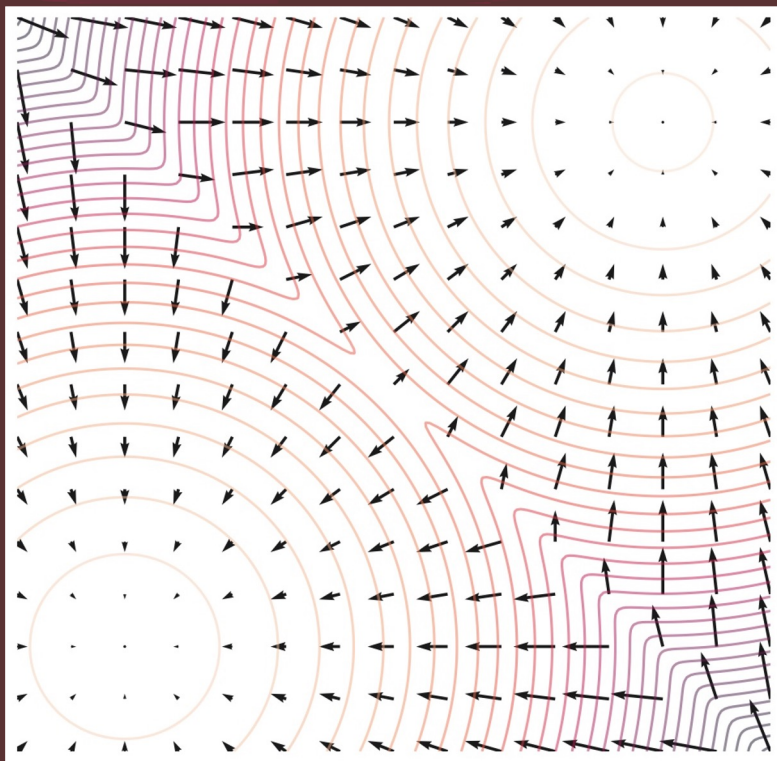
Learning to generate by iterative denoising.

# Diffusion models

- Forward diffusion (noising)
  - $x_0 \rightarrow x_1 \rightarrow \cdots x_T$
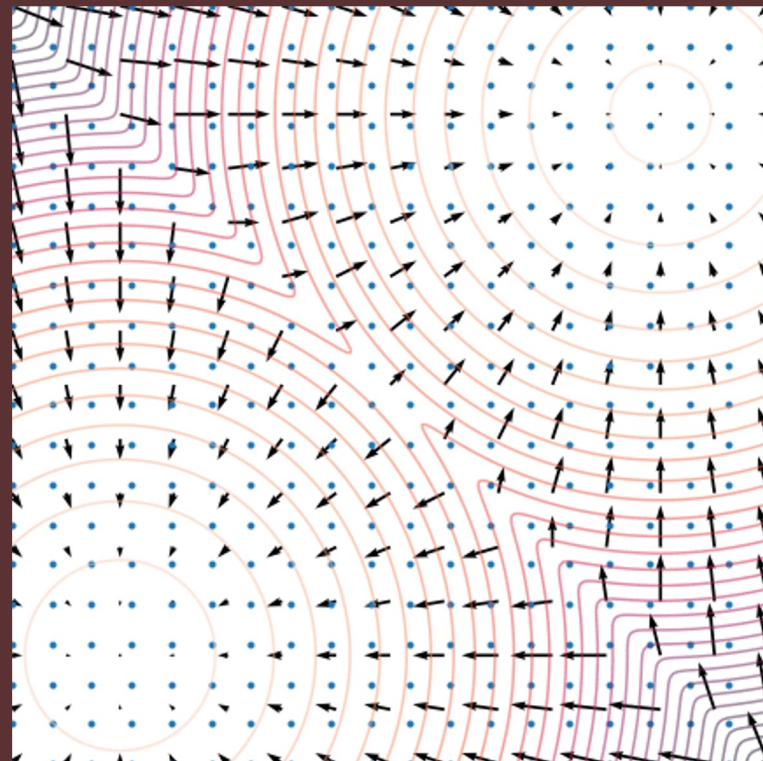  - Take a data distribution $x_0 \sim p(x)$, turn it into noise by diffusion $x_T \sim \mathcal{N}(0, \sigma^2 I)$



$\boldsymbol{x_0}$      $\boldsymbol{x_1}$                          $\boldsymbol{x_{T-1}}$      $\boldsymbol{x_T}$

- Reverse diffusion (denoising)
  - $x_T \rightarrow x_{T-1} \rightarrow \cdots x_0$
  - Sample from the noise distribution $x_T \sim \mathcal{N}(0, \sigma^2 I)$, reverse the diffusion process to generate data $x_0 \sim p(x)$

# Animation for the Reverse Diffusion



Score Vector Field



Reverse Diffusion guided by the score vector field

# Training diffusion model =
# Learning to denoise

- If we can learn a score model
$$f_\theta(x, t) \approx \nabla \log p(x, t)$$
  - Then we can denoise samples, by running the reverse diffusion equation. $x_t \rightarrow x_{t-1}$

- Score model $f_\theta: \mathcal{X} \times [0,1] \rightarrow \mathcal{X}$
  - A time dependent vector field over $x$ space.

- Training objective: Infer noise from a noised sample
$$x \sim p(x), \epsilon \sim \mathcal{N}(0, I), t \in [0,1]$$
$$\min \left\| \epsilon + f_\theta(x + \sigma^t \epsilon, t) \right\|_2^2$$
  - Add Gaussian noise $\epsilon$ to an image $x$ with scale $\sigma^t$, learn to infer the noise $\sigma$.

# Conditional denoising

- Infer noise from a noised sample, based on a condition $y$
  - $x, y \sim p(x, y), \epsilon \sim \mathcal{N}(0, I), t \in [0,1]$
  - $\min \left\| \epsilon - f_\theta(x + \sigma^t \epsilon, y, t) \right\|_2^2$

- Conditional score model $f_\theta \colon \mathcal{X} \times \mathcal{Y} \times [0,1] \to \mathcal{X}$
  - Use Unet as to model image to image mapping
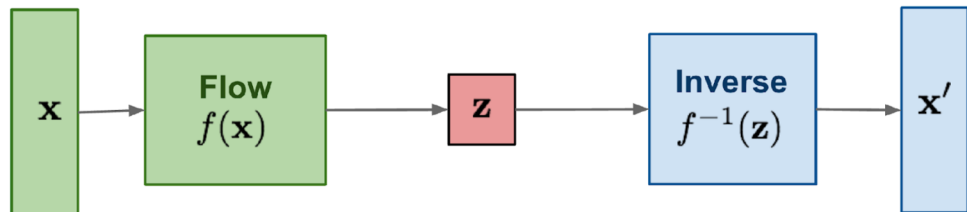  - Modulate the Unet with condition (text prompt).

# Comparing Generative Models



**GAN:** Adversarial training — Discriminator $D(\mathbf{x})$ → 0/1 → Generator $G(\mathbf{z})$

**VAE:** maximize variational lower bound — Encoder $q_\phi(\mathbf{z}|\mathbf{x})$ → $\mathbf{z}$ → Decoder $p_\theta(\mathbf{x}|\mathbf{z})$

**Flow-based models:** Invertible transform of distributions — Flow $f(\mathbf{x})$ → $\mathbf{z}$ → Inverse $f^{-1}(\mathbf{z})$

**Diffusion models:** Gradually add Gaussian noise and then reverse — $\mathbf{x}_0$, $\mathbf{x}_1$, $\mathbf{x}_2$ ... ... $\mathbf{z}$

https://lilianweng.github.io/posts/2021-07-11-diffusion-

# Modelling Score function over Image Domain

Introducing UNet



Denoising U-Net $\epsilon_\theta$

# Convolutional Neural Network



Features of **larger** scale (larger RF)
**Higher** abstraction level.

224 x 224 x 3    224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512    14 x 14 x 512    7 x 7 x 512

1 x 1 x 4096    1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

VGG

- CNN parametrizes function over images

- Motivation
  - Features are translational invariant
  - Extract feature at different scale / abstraction level

- Key modules
  - Convolution
  - Downsamping (Max-pool)

# UNet: a natural architecture for image-to-image function

# *Note:* Add Time Dependency

- The score function is *time-dependent.*
  - Target: $s(x, t) = \nabla_x \log p(x, t)$

- Add time dependency

  - Assume time dependency is spatially homogeneous.

  - Add one scalar value per channel $f(t)$

  - Parametrize $f(t)$ by MLP / linear of Fourier basis.

# How to understand prompts?

Language / Multimodal Transformer, CLIP!

# Word as Vectors: Language Model 101

- Unlike pixel, meaning of word are not explicitly in the characters.

- Word can be represented as index in dictionary
  - But index is also meaning less.

- Represent words in a vector space
  - Vector geometry => semantic relation.

**Words in a sentence**  I  love  cats  and  dogs  .

**Token Index**  328,  793,  3989,  537,  3255,  269

**Word Vectors**

# Word Vector in Context: RNN / Transformers

*N layers ……*

- Meaning of word depends on context, not always the same.

  - "I book a ticket to buy that book."

  - Word vectors should depend on context.

- Transformers let each word "absorb" influence from other words to be "contextualized"

**Transformer Block**

**Transformer Block**

I  love  cats  and  dogs  .



Transformer block

Output Embeddings

Residue, Norm

FFN

Self Attention

Residue, Norm

Input Embeddings

*More on attention later...*

# Learning Word Vectors: GPT & BERT & CLIP

- Self-supervised learning of word representation

    - Predicting missing / next words in a sentence. (BERT, GPT)

    - Contrastive Learning, matching image and text. (CLIP)

Downstream Classifier can decode:
Part of speech, Sentiment, …

# Joint Representation for Vision and Language : CLIP



- Learn a joint encoding space for text caption and image

- Maximize representation similarity between an image and its caption.

- Minimize other pairs

CLIP paper 2021

# Choice of text encoding

- Encoder in Stable Diffusion: pre-trained CLIP ViT-L/14 text encoder

- Word vector can be randomly initialized and learned online.

- Representing other conditional signals

  - Object categories (e.g. Shark, Trout, etc.):
    - 1 vector per class

  - Face attributes (e.g. {female, blonde hair, with glasses, … }, {male, short hair, dark skin}):
    - set of vectors, 1 vector per attributes

- Time to be creative!!

# Origin of Attention:
# Machine Translation (Seq2Seq)

**Original sentence**    I  love  cats  and  dogs  .    **French Translation** J'adore  les  chats  et les chiens.

**Encoder hidden state (Word Vectors)**    $e_1$  $e_2$  $e_3$  $e_4$  $e_5$  $e_6$    **Decoder hidden state (Word Vectors)**  $h_1$  $h_2$  $h_3$

- Use **Attention** to retrieve useful info from a batch of vectors.

# From Dictionary to Attention
## Dictionary: Hard-indexing

- `dic = {1 : v_1, 2 : v_2, 3 : v_3}`
  - Keys  1,2,3
  - Values  $v_1, v_2, v_3$

- `dic[2]`
  - Query 2
  - Find 2 in keys
  - Get corresponding value.

- Retrieving values as matrix vector product
  - One hot vector over the keys
  - Matrix vector product

$$
\begin{array}{ccc} \mathbf{1} & \mathbf{2} & \mathbf{3} \end{array}
$$

$$
\begin{bmatrix} & & \\ v_1 & v_2 & v_3 \\ & & \end{bmatrix} \times \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} \\ v_2 \\ \end{bmatrix}
$$

# From Dictionary to Attention
## Attention: Soft-indexing

- Soft indexing
  - Define an attention distribution $a$ over the keys

  - Matrix vector product.

  - Distribution based on similarity of query and key.

# QKV attention

- **Q**uery : what I need (*J'adore* : "I want subject pronoun & verb")
- **K**ey : what the target provide (*I* : "Here is the subject")
- **V**alue : the information to be retrieved (latent related to *Je* or *J'* )

<br>

- Linear projection of "word vector"
  - Query   $q_i = W_q h_i$
  - Key   $k_j = W_k e_j$
  - Value  $v_j = W_v e_j$

  - *$e_j$ hidden state of encoder (English, source)*
  - *$h_i$ hidden state of decoder (French, target)*

# Attention mechanism

- Compute the inner product (similarity) of key $k$ and query $q$

- SoftMax the normalized score as attention distribution.

$$a_{ij} = \text{SoftMax}\left(\frac{k_j^T q_i}{\sqrt{len(q)}}\right), \sum_j a_{ij} = 1$$

- Use attention distribution to weighted average values $v$.

$$c_i = \sum_j a_{ij} v_j$$

# Cross & Self Attention

- Cross Attention
  - Tokens in one language pay attention to tokens in **another.**

- Self Attention $(e_i = h_i)$
  - Tokens in a language pay attention to **each other.**

**French Translation** J'adore les chats

**Decoder hidden state (Word Vectors)** $h_1$ $h_2$ $h_3$

"A robot must obey the order given *it*."

"A robot must obey the order given *it*."

Query #9

it

0.1% value #4

0.1% value #3

50% value #2

30% value #1

a

robot

must

obey

| Word | Value vector | Score | Value X Score |
|:---:|:---:|:---:|:---:|
| <s> | ⬜⬜⬜ | 0.001 | ⬜⬜⬜ |
| a | ⬜⬜⬜ | 0.3 | ⬜⬜⬜ |
| robot | ⬜⬜⬜ | 0.5 | ⬜⬜⬜ |
| must | ⬜⬜⬜ | 0.002 | ⬜⬜⬜ |
| obey | ⬜⬜⬜ | 0.001 | ⬜⬜⬜ |
| the | ⬜⬜⬜ | 0.0003 | ⬜⬜⬜ |
| orders | ⬜⬜⬜ | 0.005 | ⬜⬜⬜ |
| given | ⬜⬜⬜ | 0.002 | ⬜⬜⬜ |
| it | ⬜⬜⬜ | 0.19 | ⬜⬜⬜ |
| | | | |
| | | Sum: | 🟥🟥🟥 |

# Text2Image as translation

**Source language: Words**

**Target language: Images**

Sequence Dimensions

**Encoded Word Vectors**

" A ballerina chasing her cat running on the grass in the style of Monet "

Spatial Dimensions

Channel Dimensions

**Latent State of Image**

*Patch Vectors!*

# Text2Image as translation

**Encoded Word Vectors**

Sequence Dimensions

**Latent State of Image**

Spatial Dimensions

Channel Dimensions

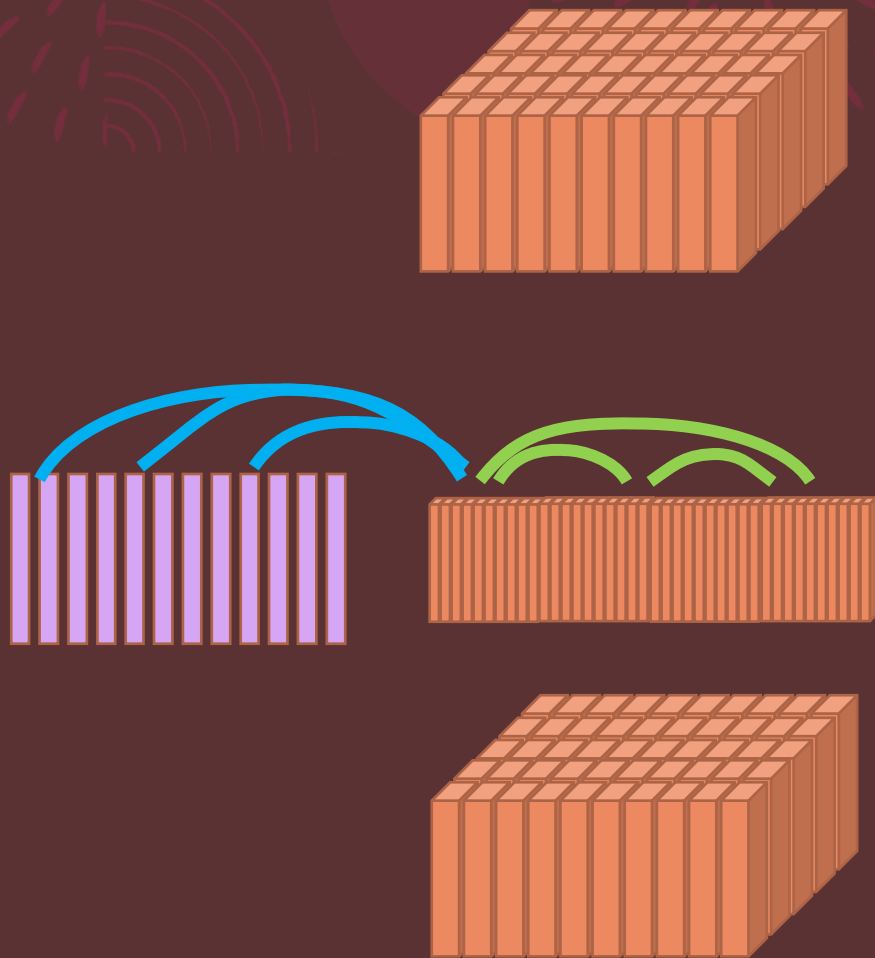" A ballerina chasing her cat running on the grass in the style of Monet "

**Cross Attention:**
Image to Words

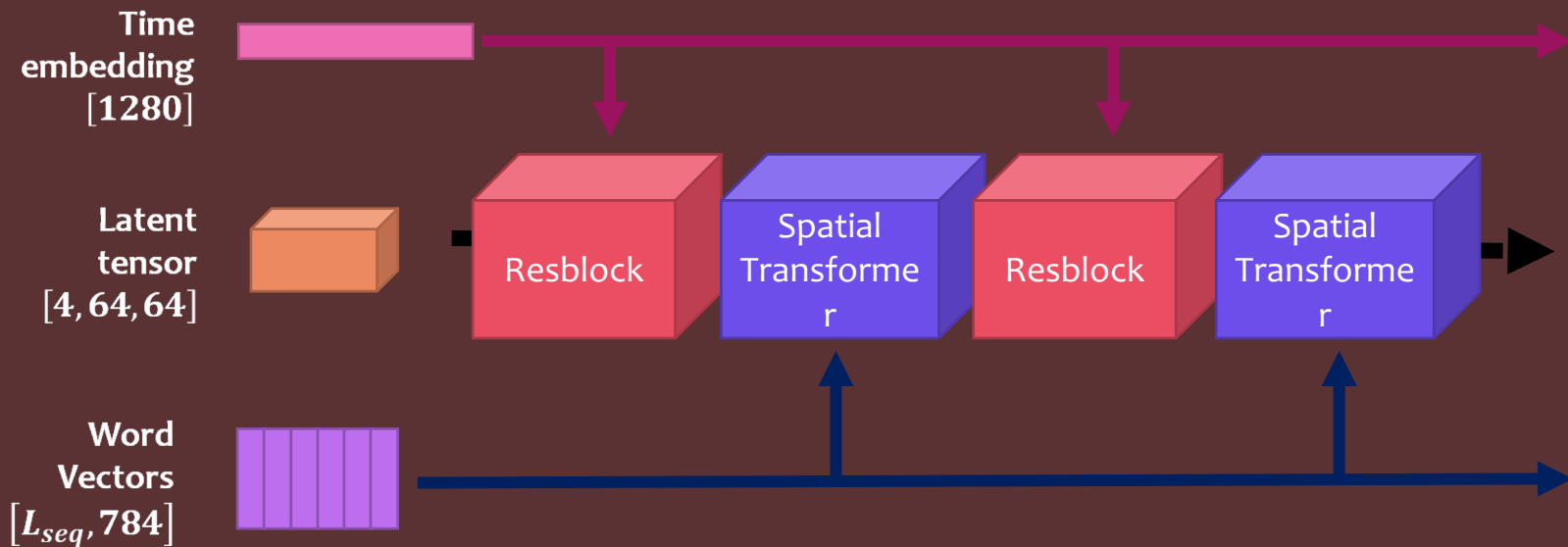**Self Attention:**
Image to Image

# Spatial Transformer

- Rearrange spatial tensor to sequence.

- Cross Attention

- Self Attention

- FFN

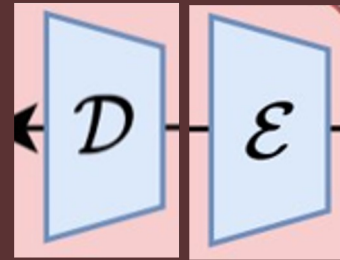- Rearrange back to spatial tensor (same shape)

# Spatial transformer + ResBlock (Conv layer)



- Alternating Time and Word Modulation
- Alternating Local and Nonlocal operation

# Diffusion in Latent Space

Adding in AutoEncoder

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models, *CVPR*

# Diffusion in latent space

- Motivation:
  - Natural images are high dimensional
  - but have many redundant details that could be compressed / statistically filled out

- Division of labor
  - Diffusion model -> Generate low resolution sketch
  - AutoEncoder -> Fill out high resolution details

- Train a VAE model to compress images into latent space.
  - $x \rightarrow z \rightarrow x$
- Train diffusion models in latent space of $z$.
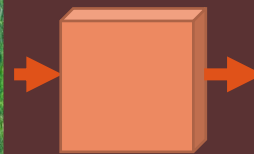
DownSampling

**32 pix**          **180 pix**

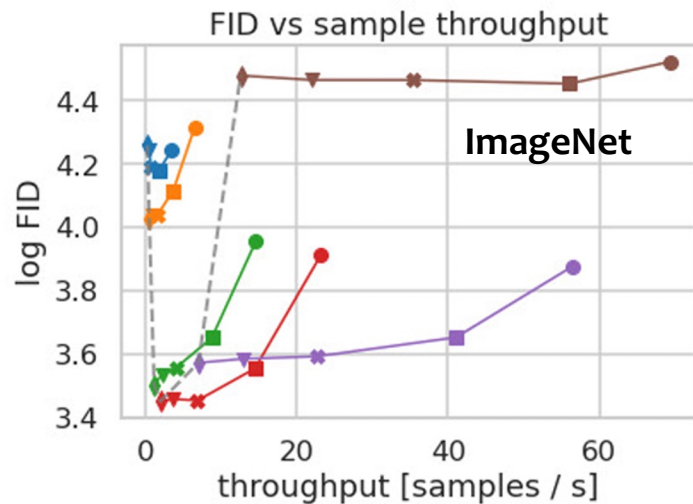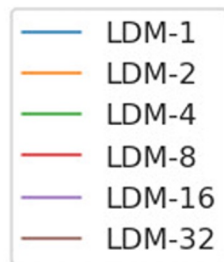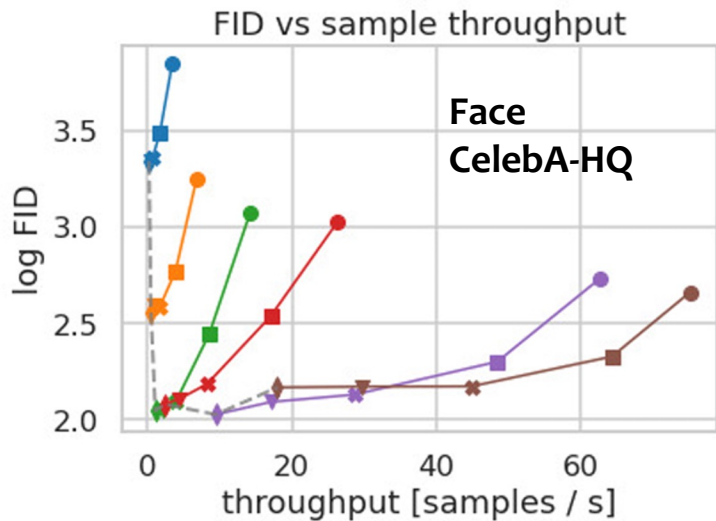$d = 2352$          $d = 97200$

$x$
$[3,512,512]$

$z$
$[4,512/f, 512/f]$

$\hat{x}$
$[3,512,512]$

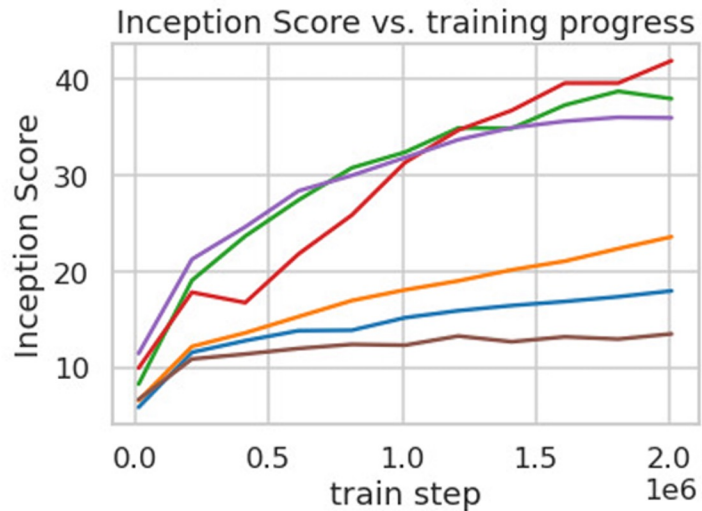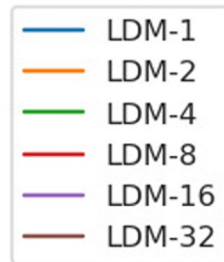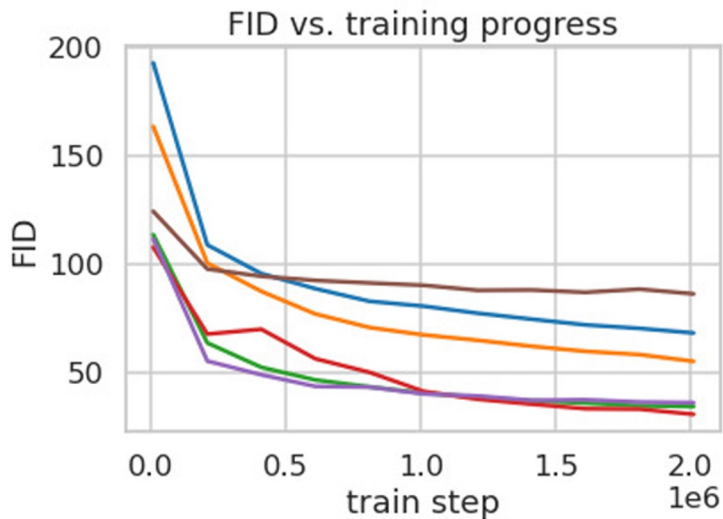# Spatial Compression Tradeoff

- LDM-$\{f\}$. $f$ = Spatial downsampling factor
  - Higher $f$ leads to faster sampling, with degraded image quality (FID $\uparrow$)
  - Fewer sampling steps leads to faster sampling, with lower quality (FID $\uparrow$)

# Spatial Compression Tradeoff

- LDM-$\{f\}$. $f$ = Spatial downsampling factor
  - Too little compression $f = 1,2$ or too much compression $f = 32$, makes diffusion hard to train.

# Future Direction

- Introduce more modal, such as video-text, music-text, music-image.

- Speed up the generation

# Thank you