

Stable Diffusion.

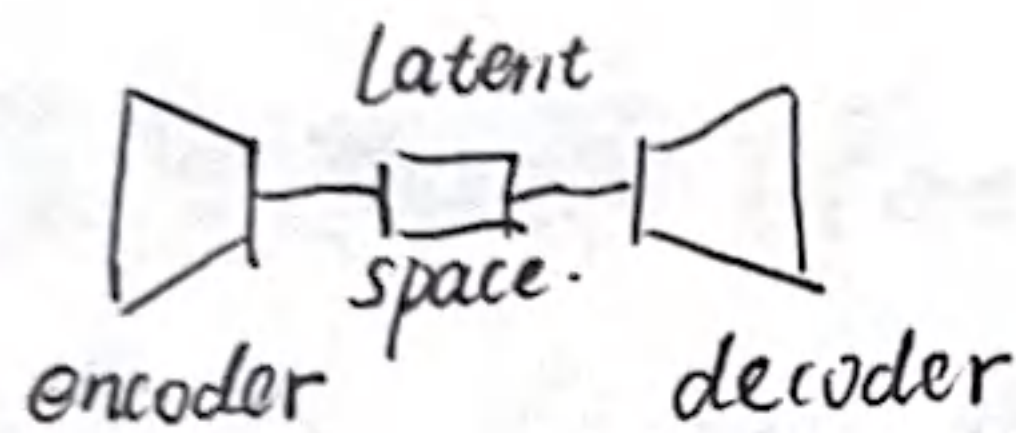
1. 使用了 CLIP (Contrastive Language-Image Pre-training) 中的 Text Encoder

2. Variational Auto Encoder.

(1) 用于压缩前/逆向过程的维度以提高计算效率
→ latent diffusion.

→ 将学习 data distribution \Rightarrow latent representation of the data.

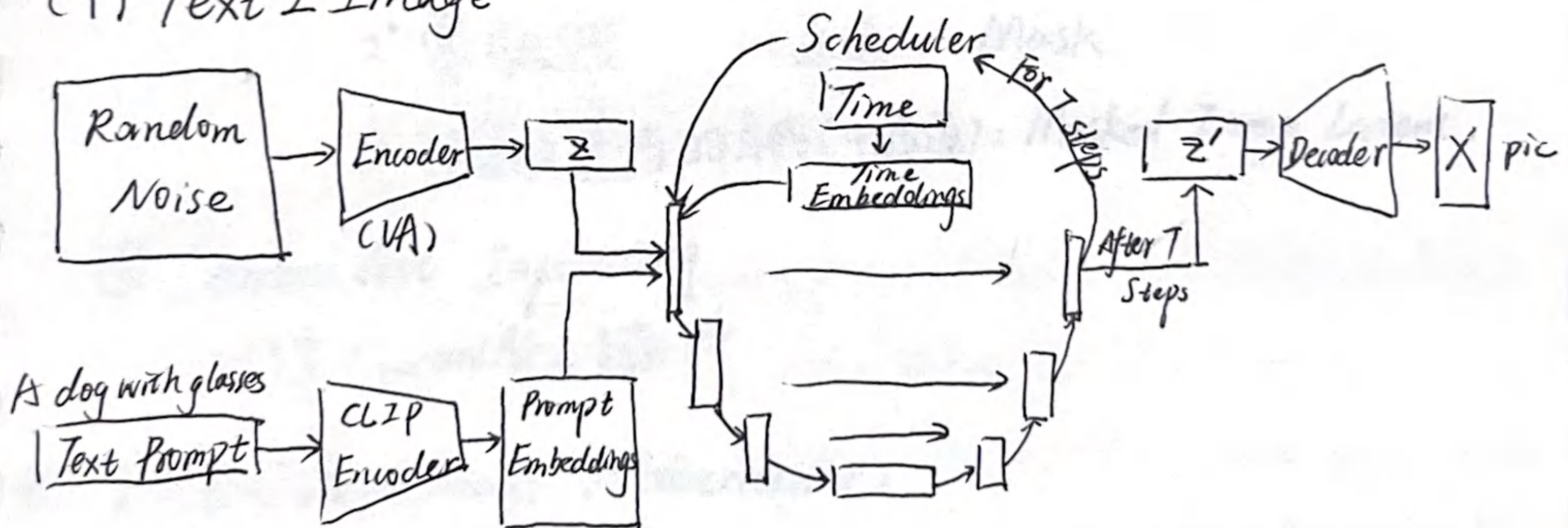
(2) 不压缩, 同时也学习了一个 latent space.



↑ represents the parameters of a multivariate distribution (多元分布)
(Gaussian) 学的是 μ 和 σ .

3. Architecture.

(1) Text 2 Image



(2) Image 2 Image

将 Random Noise 换成真实图片; encoder 得到的 z 是 z_0 . Add noise to Latent.

- 添加的噪声少, 生成的自由度 \uparrow radically
- 少的噪声 means 不能从根本上改变图片.
- 噪声多少 \Leftrightarrow 对原图有多大注意力

(3) In-Painting. (fool the model) \Rightarrow 引入约束

① 推理时干预 (Inference-time Strategy)

• 代表作: RePaint

• 思想: 不重新训练模型, 在去噪过程作弊.

$$X_{t-1}^{final} = \underset{\substack{\uparrow \\ \text{被遮住的部分} \\ \text{使用模型推理结果}}}{Mask \cdot X_{t-1}^{pred}} + (1 - Mask) \cdot \underset{\substack{\uparrow \\ \text{保留区域: 不是直接将原图 } X_0 \text{ 贴上去} \\ \text{而是前向到 } t-1 \text{ 强度的噪声}}}{X_{t-1}^{known}}$$

• 问题: 边缘不连贯.

\Rightarrow 生成 X_{t+1} 后, 强行加噪回 X_t , 反复横跳.

计算量大

② 模型微调 (Fine-tuning Strategy)

• 代表作: Stable Diffusion Inpainting Model.

• 思想: 将输入通道 $4 \rightarrow 9 \Rightarrow$ Unet.

1° 当前噪声图 (4通道): ϵ_t

2° 蒙片版图 (1通道): Mask

3° 被蒙片遮挡的原图 (4通道): Masked Image Latent.

③ ControlNet Inpainting

训练 ControlNet 插件.

4. VAE (Variational Autoencoder)

(1) • 原理: Not learning how to compress the data but learning a latent space which are the parameters of a multivariate Gaussian distribution.

\rightarrow Actually, the VAE is trained to learning the μ and σ .

\rightarrow Then we will sample the distribution.

\downarrow
 $\log \sigma$.

5. 核心参数

① strength: 重绘程度

- 决定了生成图和原图有多像 图片还未完全变成噪声
- 代码逻辑: 起始噪声水平 (从中间步开始降噪)

② do_cfg = True: 启用分类器自由引导.

- 是否听懂人话
- 代码逻辑: 如果为 True, 模型的每步预测会同时计算 2 个结果:
 - 1° 有提示词预测 (output-cond)
 - 2° 无 ~ (output-uncond): 完全盲猜, 根据负向 prompt 或空内容去预测.
 - 3° 将 2 种 output 结合
- 如果不开启, 生成图虽然自然, 但不会理会 prompt.

③ cfg-scale: 引导系数

- 代码逻辑: 最终结果 = 无引导结果 + 系数 (有引导 - 无引导)
- < 2 : 不理睬 prompt
- > 15 : 强行匹配每一个词.

6. 代码小思考

(1) 调度策略: 原 DDPM 用的 Linear Schedule, SD 中用的无开方是 做插值再平方回去的 Scaled Linear Schedule. why?

Ans. DDPM 是在像素空间训练的, SD 是在 Latent space 上训练的, 对于方差和分布更敏感. 采用 SLS, 让 β 在初期增长得慢一些, 适合 Latent Space 的信噪比变化规律.

(2) 项目中采用的是 Respaced DDPM, 而非 DDIM?

Ans. 因为每步噪声采样的 Variance $\neq 0$, 如果为 0 的话才是 DDIM.

(3) VAE 中添加的噪声和 diffusion 前添加的噪声的区别?

Ans. VAE 把图片压缩后, 得到的是一个高维分布 (μ, σ) ,

需要从中取出一个具体的 Latent 向量 z 为后面的流程用，
 因此需要重参数化采样: $Latent = Mean + Variance \times \frac{Noise}{\sigma}$.
 没有这个噪音，VAE 无法训练退化为普通的编码器。

$$(z | x) \sim N(\mu(x), \sigma^2(x)) \Rightarrow \frac{p(x|z)p(z)}{p(x)} = p(z)$$

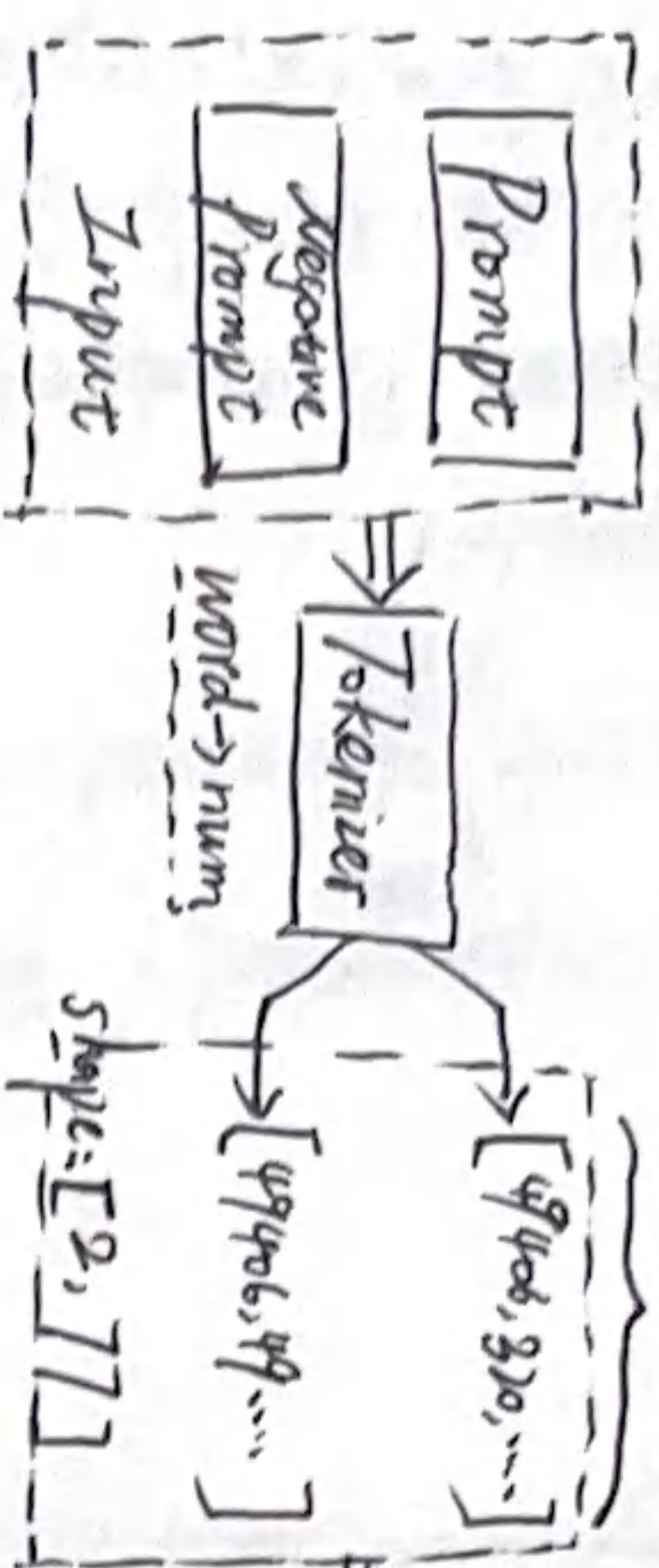
$$\begin{aligned} 3 \sqrt{\sigma^2(x)} + 0.5 \sqrt{6 \ln m + k} &= 3\sigma + (-\ln m + 0.5 \ln k) = 2\sigma \\ 3 \sqrt{\sigma^2(x)} + 0.5 \sqrt{6 \ln m} &= \end{aligned}$$

$$(z | x) \sim N(\mu(x), \sigma^2(x)) \Rightarrow$$

$$(z | x) \sim N(\mu(x), \sigma^2(x)) \Rightarrow \mu(x) = \mu(x) + \sigma(x) = \mu$$

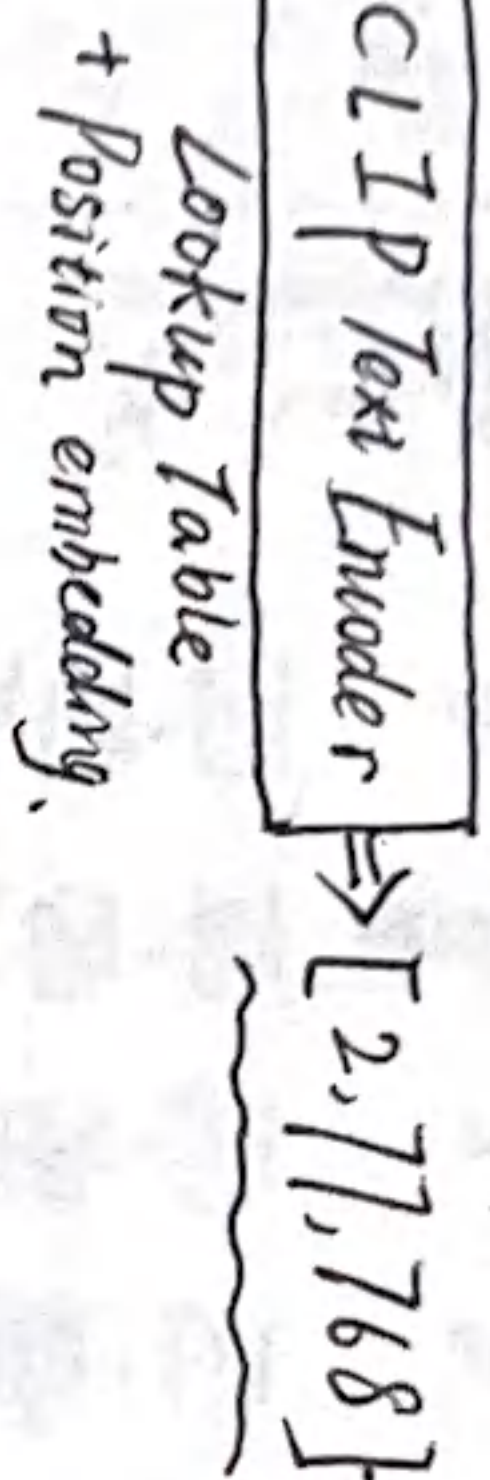
$$\mu(x) = \mu(x) + \sigma(x) = \mu$$

Step 1: Conditioning

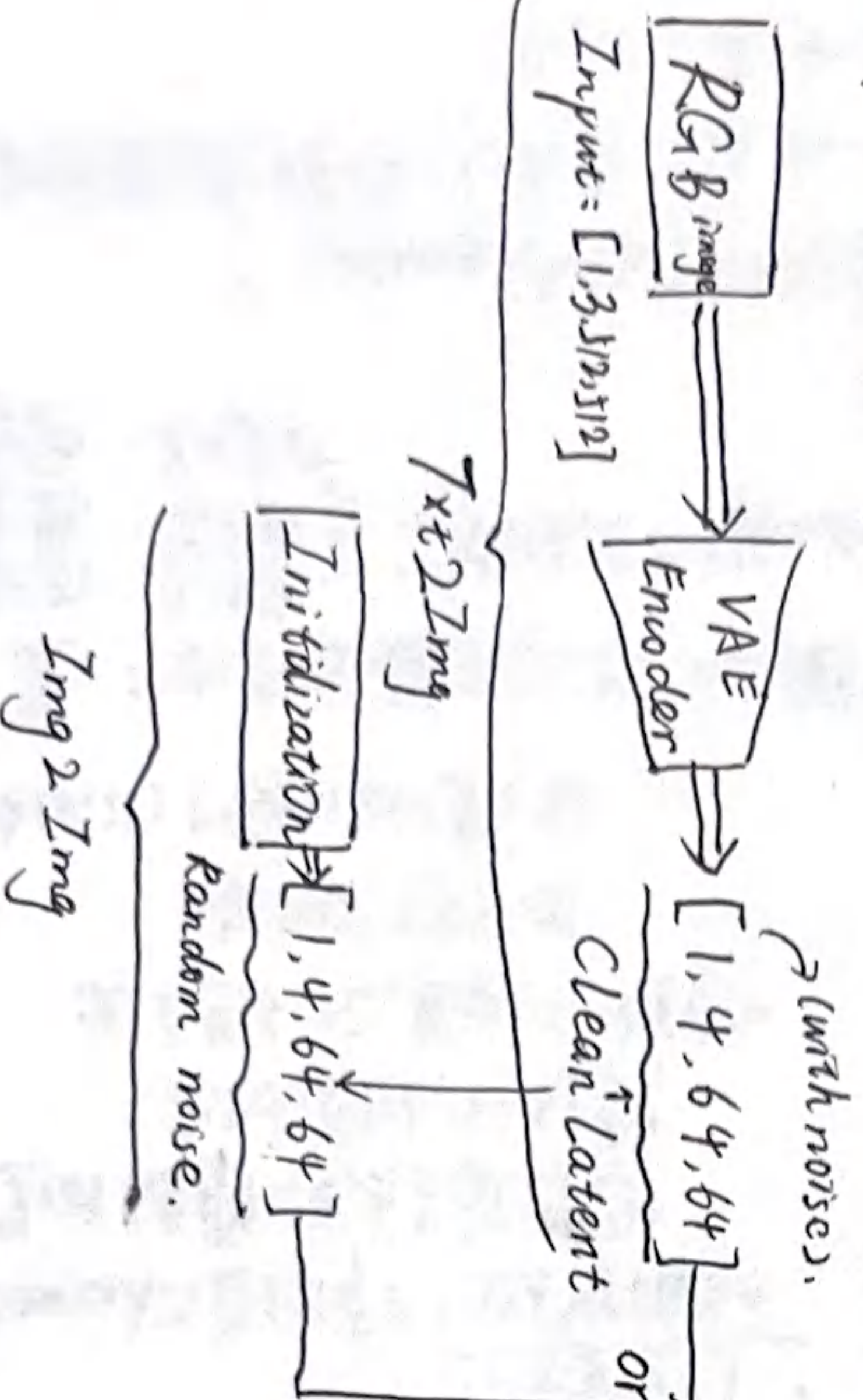


length = 77

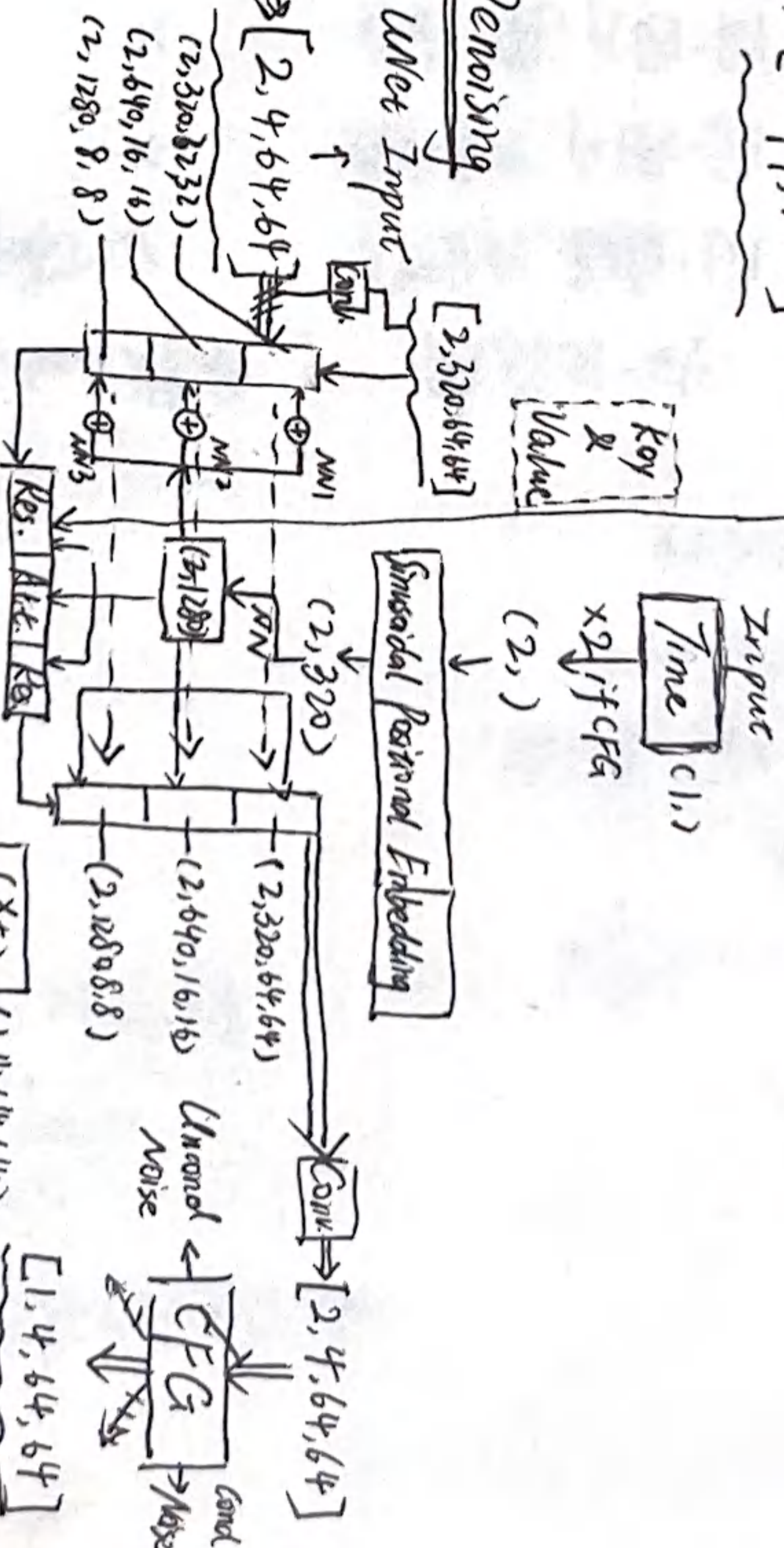
Cross-Attention = Key, Value.



Step 2: Latent Preparation

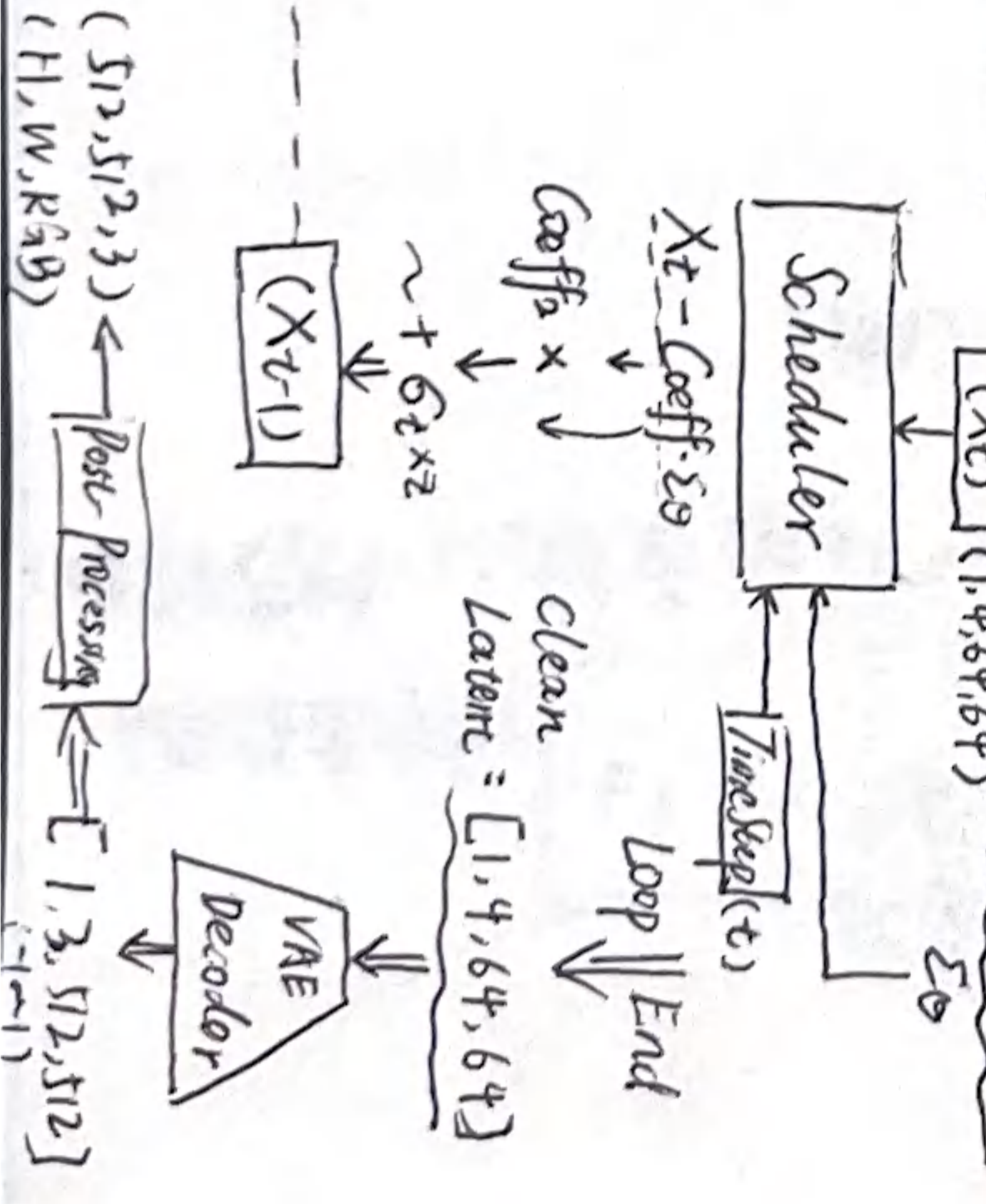


Step 3: Denoising



Step 4: Decoding

LOOP x 50



代码知识:

1. Silu



$$\text{Silu}(x) = x \cdot \sigma(x) = \frac{x}{1 + e^{-x}}$$

优势: ① 0点可导

② 负值区梯度有传递,

防止神经元死亡.

③ 隐式正则化

2. Normalization

$$y = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta \quad (\text{不改变形状})$$

标准化 \rightarrow 仿射变换 (γ, β 可学习参数)

[图]

	归一化范围	依赖 batch-size	应用场景	比喻
Batch Norm	同一通道, 所有照片	✓ (劣势)	CNN	全校单科排名
Layer Norm	同一图片, 所有通道	x	NLP	个人综合总分
Instance Norm	同一图片, 单一通道	x	风格迁移 GAN	个人单科
Group Norm	同一图片, 一组通道	x	VAE, Stable Diffusion	个人理/文特

3. 卷积输出公式

$$H_{out} = \left\lceil \frac{H_{in} + P - K}{S} \right\rceil + 1$$

• 1x1 卷积:

物理意义: 作用在 Channel 维度上的全连接层
把每个像素点的特征向量做一维矩阵乘法让特征间相互融合.

作用: ① 通道间信息交流

② 升维与降维

③ 增加非线性 (后面加激活函数).

4. nn.Sequential: 只能一种输入 x, 然后把输出传给下一层.

↓ 封装

• Switch Sequential: 处理 UNet 每层不同的需求

↓ 封装

• Module List : 文件夹

卷积层: 图像 x

残差层: 图像 x + 时间 t

注意力层: 图像 x + 文本 context.

▷ Why UNet 用 Module List 包裹许多 Switch Sequential?

Ans. UNet 的编码器 (Downloader) 每经过一层都需保存当前特征图, 以便后面解码器 (Uploader) 的每一层 concat.