

ControlNet

Links

<https://learnopencv.com/controlnet/#How-ControlNet-Works?> - Кратко

<https://arxiv.org/abs/2302.05543> - Стаття

https://llyasviel.github.io/misc/202309/cnet_supp.pdf - supplementary materials.

Проблема

The largest datasets for various specific problems (e.g., object shape/normal, human pose extraction, etc.) are usually about 100K in size, which is 50,000 times smaller than the LAION-5B [79] dataset

The direct finetuning or continued training of a large pretrained model with limited data may cause overfitting and catastrophic forgetting

Общая идея

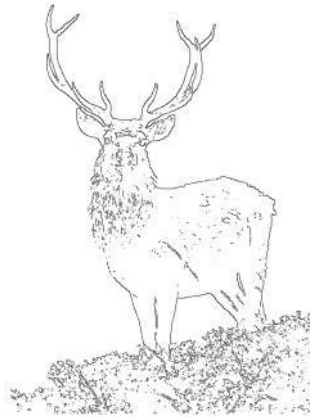
- разработать стратегию fine-tune для больших моделей не требующую огромных выч. мощ.
- как можно более четко управлять доменом

Решение:

- end-to-end архитектура для обучения моделей специфическим задачам
- скопируем веса (не все) и будем учить не их, а их копии



Source image
(for canny edge detection)

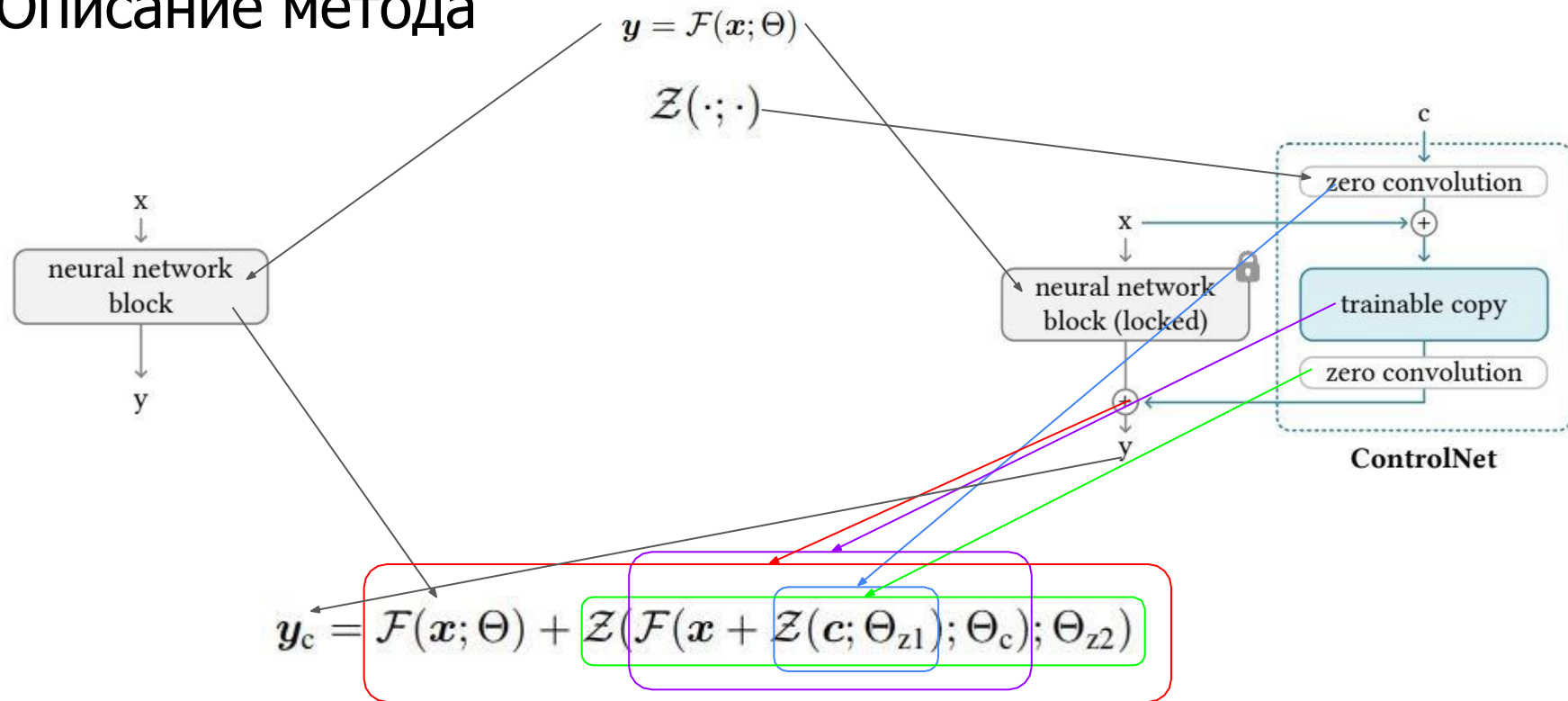


Canny edge (input)



Generated images (output)

Описание метода



where non-zero gradients are obtained and the neural network begins to learn. In this way, the zero convolutions become a unique type of connection layer that progressively grows parameters from zero to optimized values in a learned way.

Gradient Calculation of Zero Convolution Layers

$$\mathcal{Z}(I; \{W, B\})_{p,i} = B_i + \sum_j^c I_{p,j} W_{i,j}.$$

$$\begin{cases} \mathcal{Z}(c; \Theta_{z1}) = 0 \\ \mathcal{F}(x + \mathcal{Z}(c; \Theta_{z1}); \Theta_c) = \mathcal{F}(x; \Theta_c) = \mathcal{F}(x; \Theta) \\ \mathcal{Z}(\mathcal{F}(x + \mathcal{Z}(c; \Theta_{z1}); \Theta_c); \Theta_{z2}) = \mathcal{Z}(\mathcal{F}(x; \Theta_c); \Theta_{z2}) = 0 \end{cases}$$

$$y_c = y$$

$$\mathcal{Z}(I; \{W, B\})_{p,i} = B_i + \sum_j^c I_{p,i} W_{i,j}$$

$$\begin{cases} \frac{\partial \mathcal{Z}(I; \{W, B\})_{p,i}}{\partial B_i} = 1 \\ \frac{\partial \mathcal{Z}(I; \{W, B\})_{p,i}}{\partial I_{p,i}} = \sum_j^c W_{i,j} = 0 \\ \frac{\partial \mathcal{Z}(I; \{W, B\})_{p,i}}{\partial W_{i,j}} = I_{p,i} \neq 0 \end{cases}$$

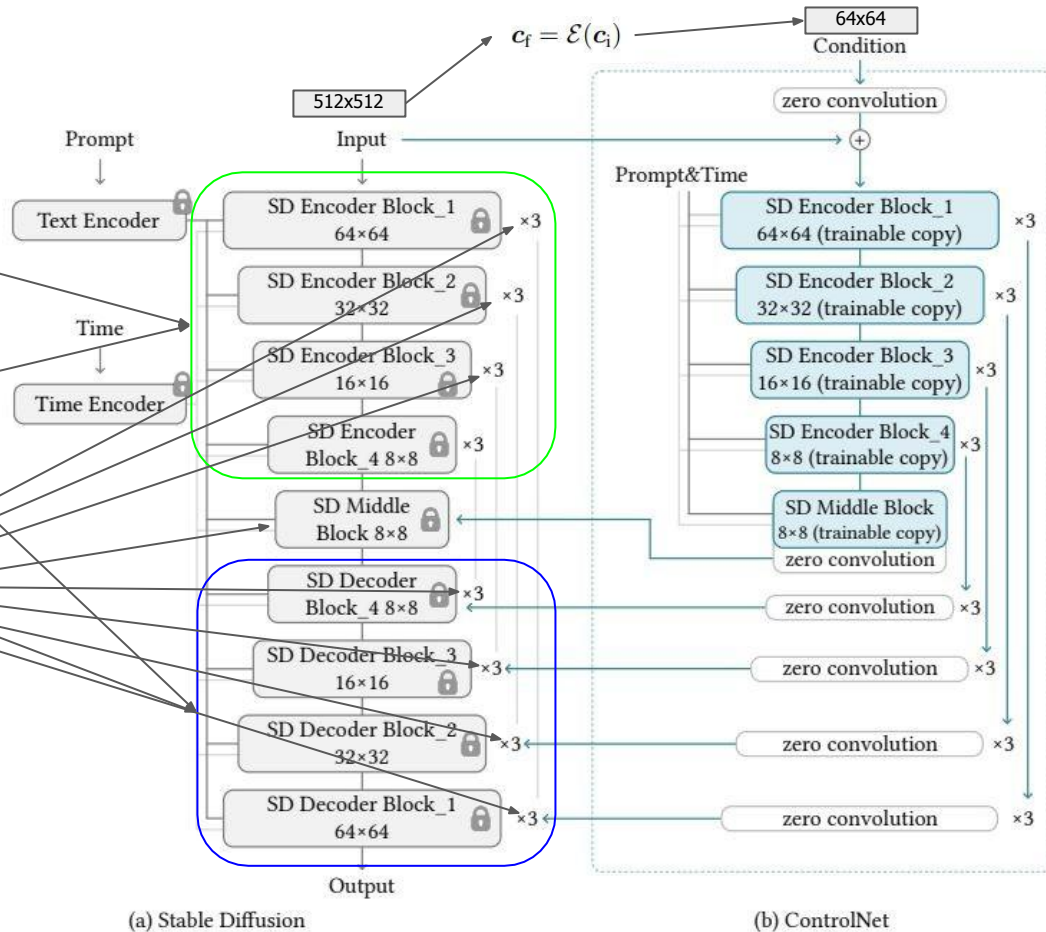
$$W^* = W - \beta_{lr} \cdot \frac{\partial \mathcal{L}}{\partial \mathcal{Z}(I; \{W, B\})} \odot \frac{\partial \mathcal{Z}(I; \{W, B\})}{\partial W} \neq 0$$

$$\frac{\partial \mathcal{Z}(I; \{W^*, B\})_{p,i}}{\partial I_{p,i}} = \sum_j^c W_{i,j}^* \neq 0$$

Hadamard
product

Control Net + SD

- Encoder: 12 blocks
- Decoder: 12 blocks
- One middle block
- 4 down-sampling
- 4 upsampling
- 17 blocks res-net + 2 ViT
- Textencoder: Clip



(a) Stable Diffusion

(b) ControlNet

Обучение

- z_0 начальное изображение
 z_t зашумленное изображение на последнем шаге
 t кол-во шагов зашумления
 c_t текстовый промт
 c_f наше условие которому учим
 ϵ_θ семейство зашумляющих сетей

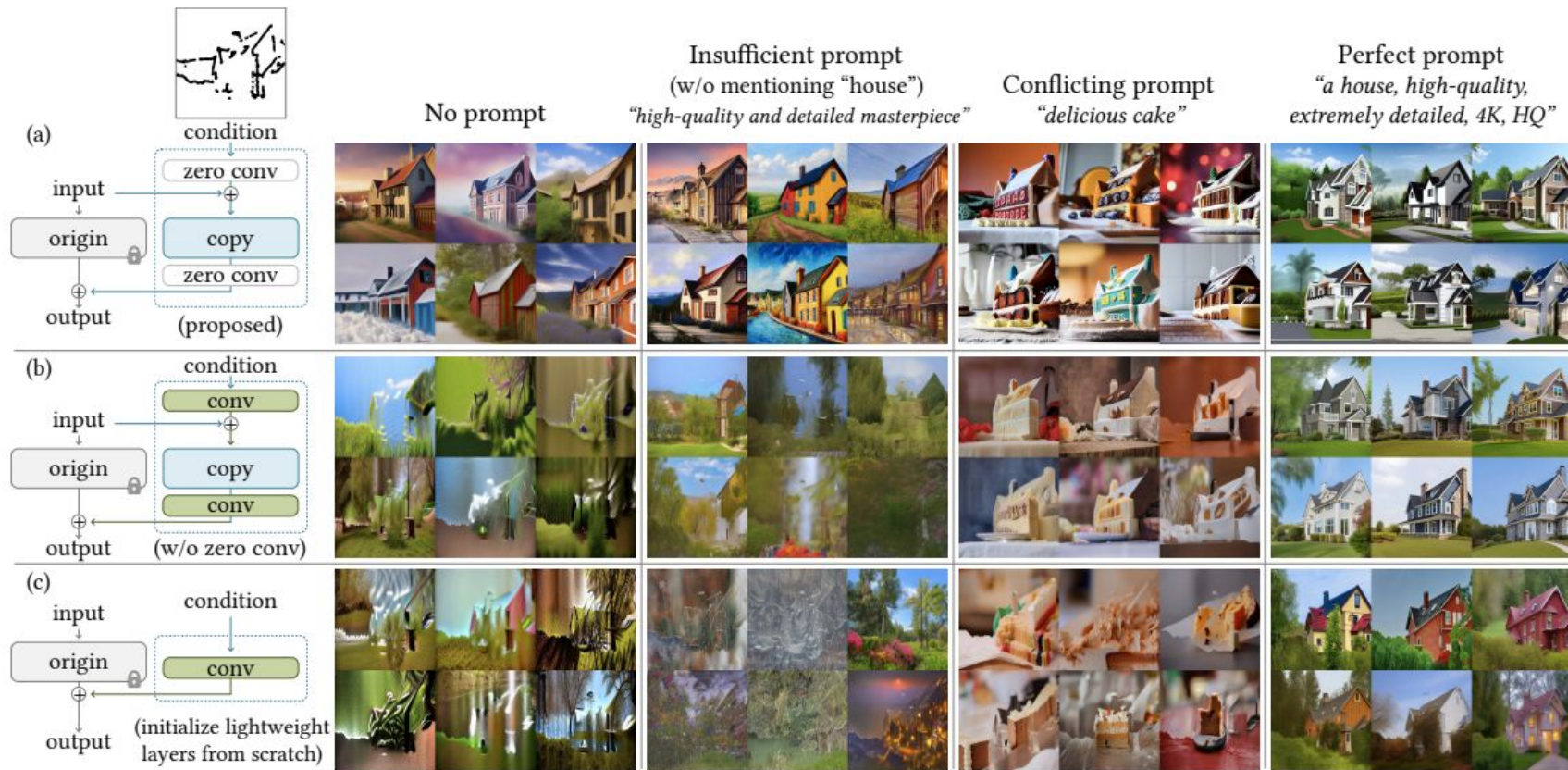
лосс функция

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

Детали обучения:

- рандомно меняем промты на пустые
- Small-Scale Training
 - “SD Middle Block” and “SD Decoder Block 1,2,3,4”
 - убрать skip-connect из D увеличит скорость в 1.6 (потом можно вернуть и дотюнить)
- Large-Scale Training
 - если много GPU
 - если много данных (пару лимонов)
 - можно учить ControlNet 50k итераций, а потом разлочить и учить всю диффузию

Ablation Study Zero Convolution



Результаты

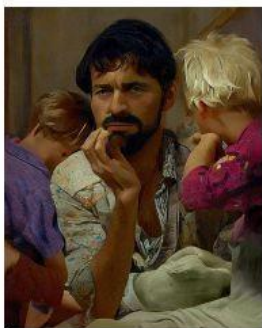
- имеем подход позволяющий сократить ресурсы при дообучении
- можно брать небольшой датасет и не оверфитнуться!
- бери доп. условие которое надо и тюнь диффузию
- на 23% нужно больше GPU и на 34% больше времени на 1 итерацию по сравнению со SD

Canny Edge

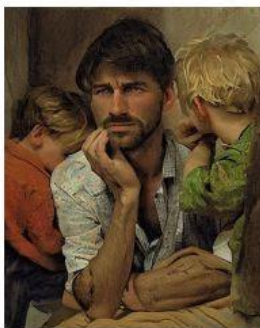
Input (Canny Edge)



Default



Automatic Prompt



“a man with beard sitting with two children”

User Prompt



“mother and two boys in a room, masterpiece, artwork”

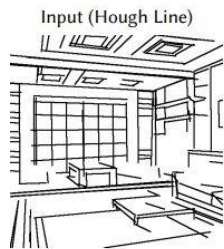


“a man in a suit and tie”



“a man in a white suit and tie”

Hough lines



Default



Automatic Prompt



"a living room with a couch and a window"



"a modern house with windows"



"a building in a city street"

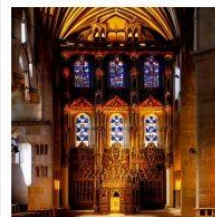
User Prompt



"a fantastic living room made of wood"



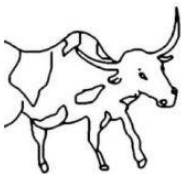
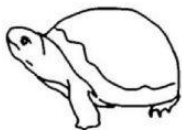
"a minecraft house"



"inside a gorgeous 19th century church"

Human scribbles

Input (User Scribble)



Default



Automatic Prompt



"a turtle in river"

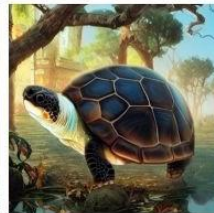


"a cow with horns standing in a field"



"a digital painting of a hot air balloon"

User Prompt



"a masterpiece of cartoon-style turtle illustration"



"a robot ox on moon, UE5 rendering, ray tracing"



"magic hot air balloon over a lit magic city at night"

Controlling Stable Diffusion with Openpose

