

# CONTROLNET

**by Klyuchnikova Ulyana**

# MOTIVATION

Let's say we have a text-to-image diffusion model but we want to add **conditional controls**

“cat crying”



“cat crying”

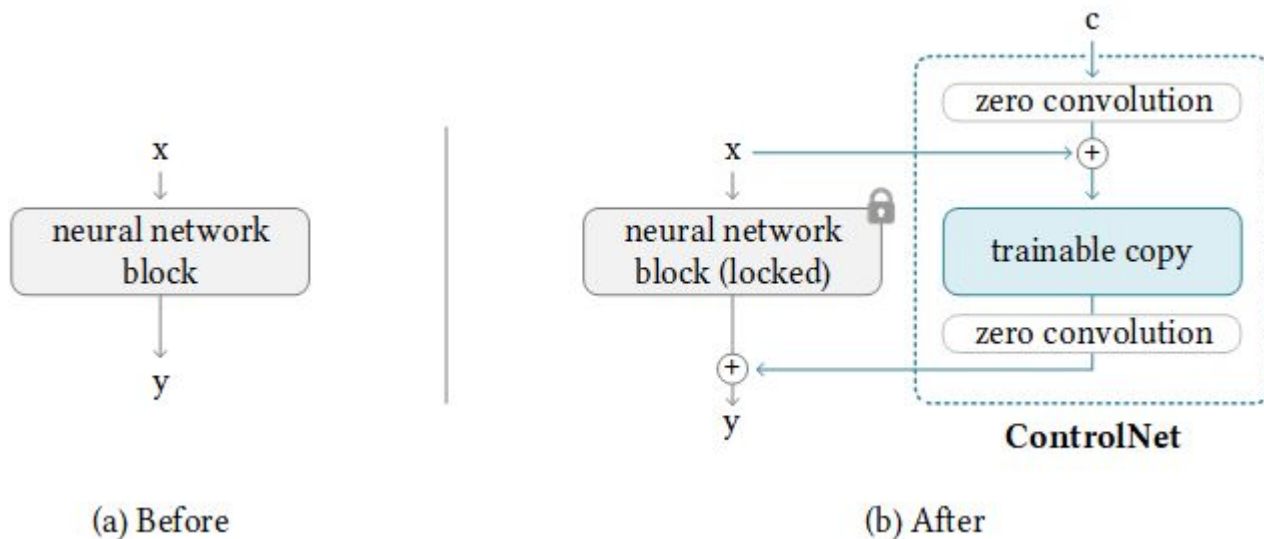
+



# GENERALLY SPEAKING...

- control the final image generation through various techniques like pose, edge detection, depth maps, etc
- end-to-end architecture
- robust on small datasets
- as fast as fine-tuning but better
- can scale to large amount of data
- does not change the initial network
- allows to train a “home-made” model

# IDEA



1. Copy model weights to locked network and trainable copy
2. Add zero convolutions – those are basically Conv 1x1 initialized by 0, so at first iteration control net has zero input

# IN FORMULAS

**x** - input like  $\mathbb{R}^{h \times w \times c}$

**y** - output

**c** - conditional vector

**$\theta$**  - original parameters

**$\theta_c$**  - trainable copy

**z** - zero 1x1 conv

Was:  $y = \mathcal{F}(x; \Theta)$

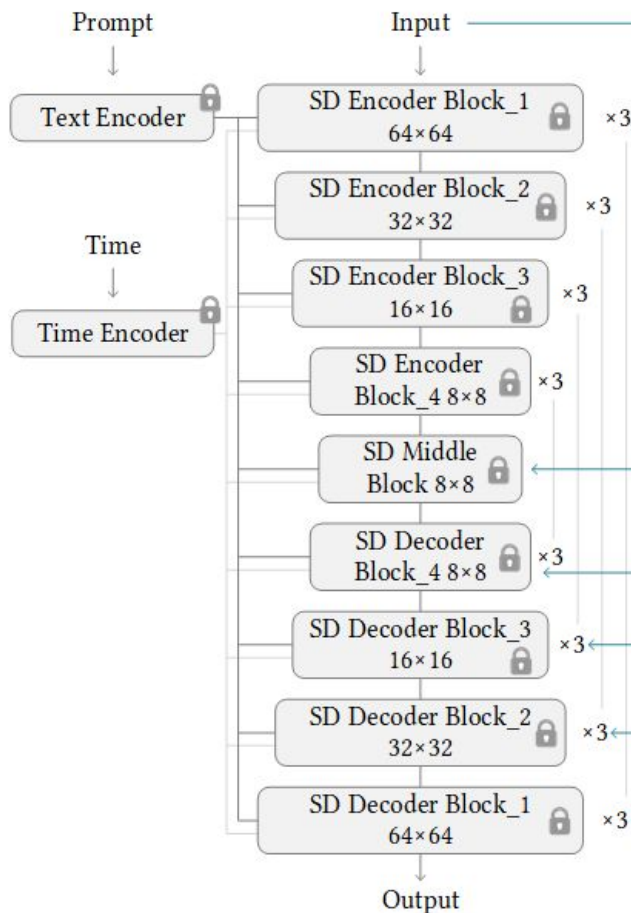
Became:  $y_c = \mathcal{F}(x; \Theta) + \mathcal{Z}(\mathcal{F}(x + \mathcal{Z}(c; \Theta_{z1}); \Theta_c); \Theta_{z2})$

# TEXT2IMAGE

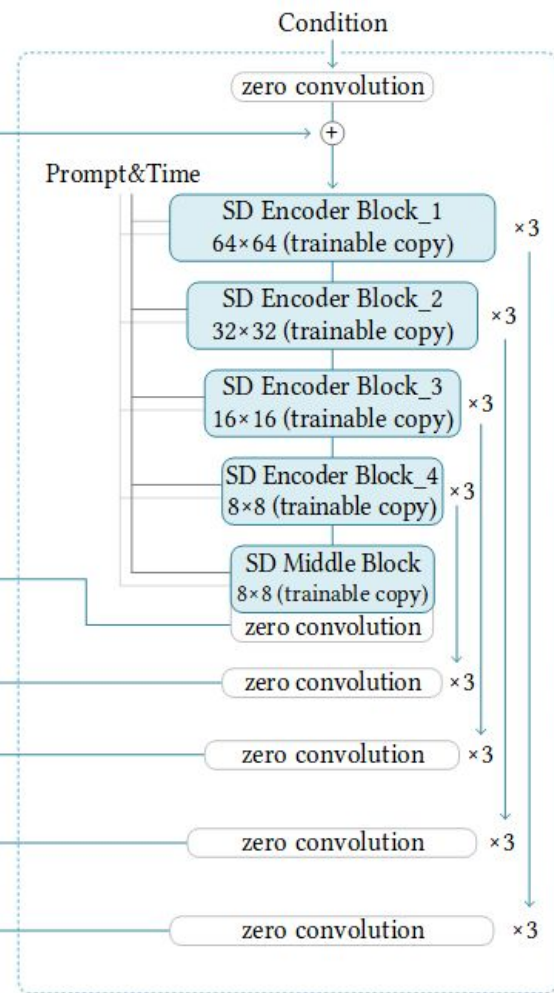
Stable Diffusion is essentially a U-Net

12 encoder blocks  
+ 1 middle  
+ 12 decoder

images as well as  
condition are  
originally 512x512  
but processed into  
64x64



(a) Stable Diffusion



(b) ControlNet

“crying cat”  
Prompt

Text Encoder

Time

Time Encoder

SD Encoder Block\_1  
64×64

×3

SD Encoder Block\_2  
32×32

×3

SD Encoder Block\_3  
16×16

×3

SD Encoder Block\_4  
8×8

×3

SD Middle  
Block 8×8

×3

SD Decoder  
Block\_4 8×8

×3

Input

512×512

$\epsilon$

Condition 64×64

zero convolution

Prompt&Time

SD Encoder Block\_1  
64×64 (trainable copy)

×3

SD Encoder Block\_2  
32×32 (trainable copy)

×3

SD Encoder Block\_3  
16×16 (trainable copy)

×3

SD Encoder Block\_4  
8×8 (trainable copy)

×3

SD Middle Block  
8×8 (trainable copy)

zero convolution

×3

# TRAINING DETAILS

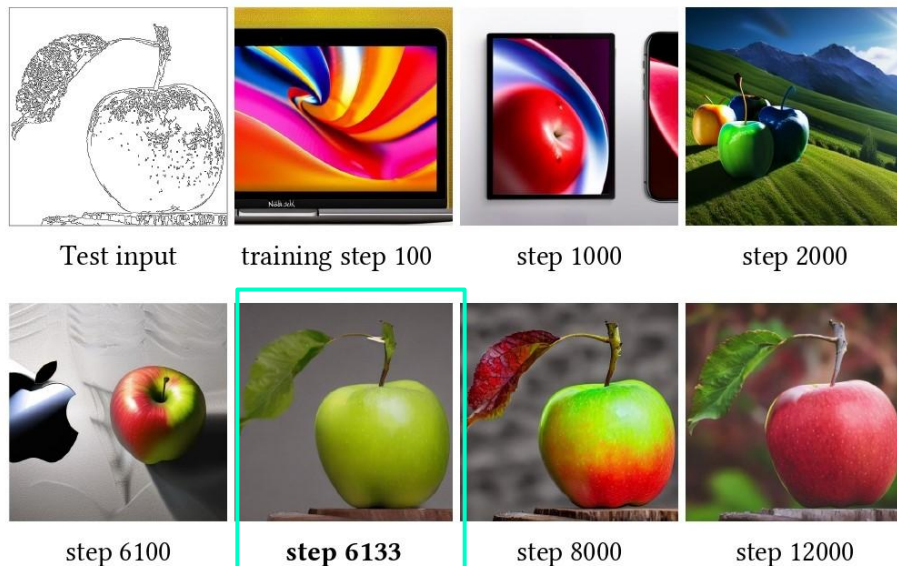
noisy image

condition

network

$$\mathcal{L} = \mathbb{E}_{\mathbf{z}_0, \mathbf{t}, \mathbf{c}_t, \mathbf{c}_f, \epsilon \sim \mathcal{N}(0,1)} \left[ \|\epsilon - \epsilon_{\theta}(\mathbf{z}_t, \mathbf{t}, \mathbf{c}_t, \mathbf{c}_f)\|_2^2 \right]$$

- 50% of prompts are replaced by empty strings to encourage control net to learn information from condition
- sudden convergence phenomenon







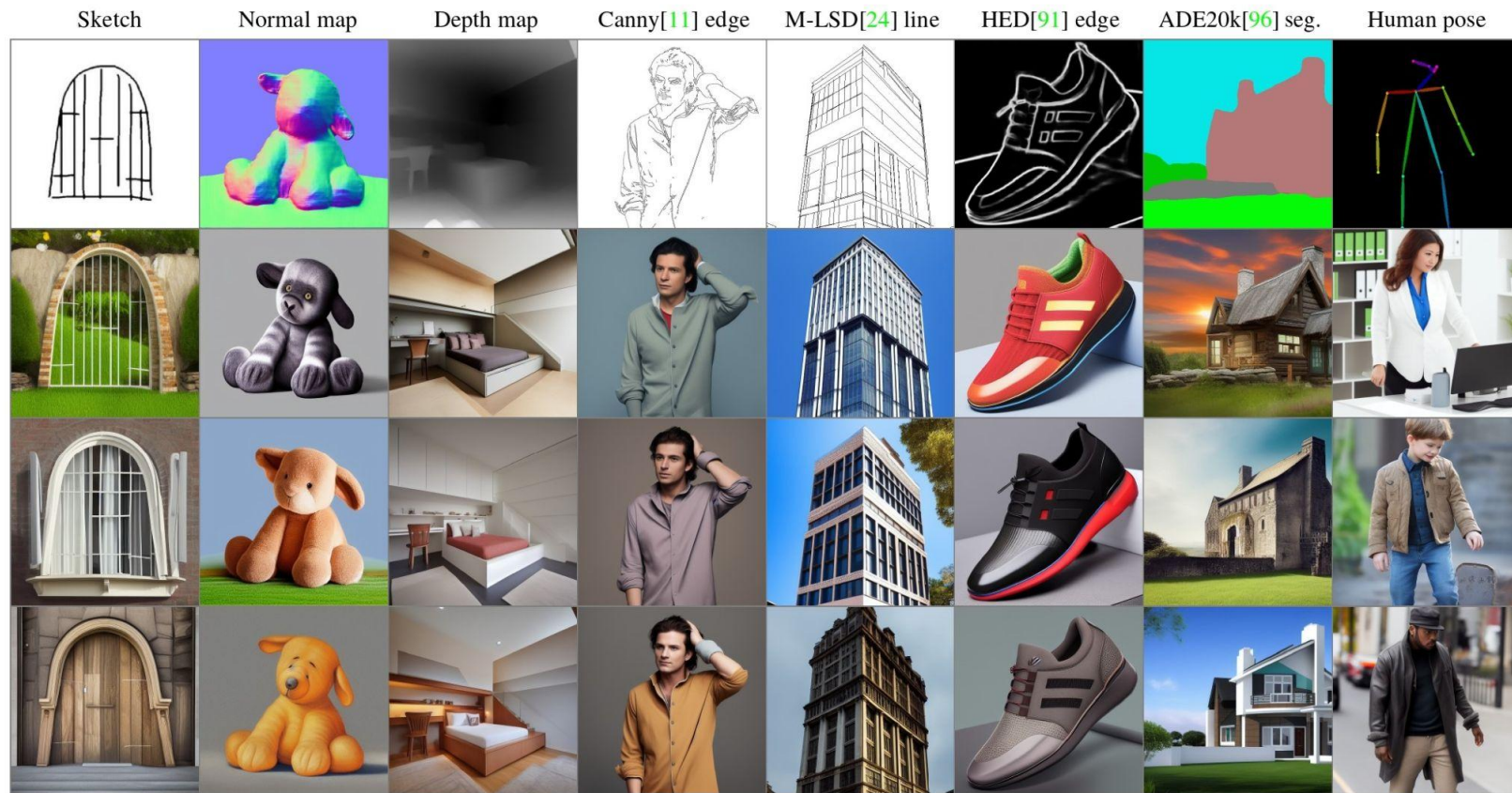
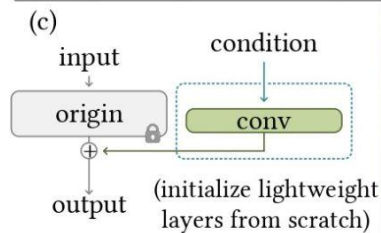
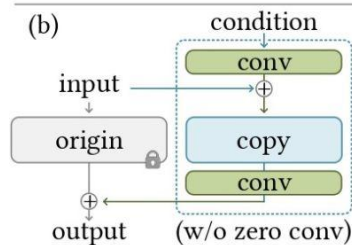
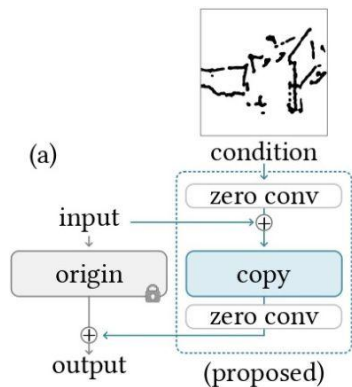


Figure 7: Controlling Stable Diffusion with various conditions **without prompts**. The top row is input conditions, while all other rows are outputs. We use the empty string as input prompts. All models are trained with general-domain data. The model has to recognize semantic contents in the input condition images to generate images.



# ABLATIVE STUDY



No prompt



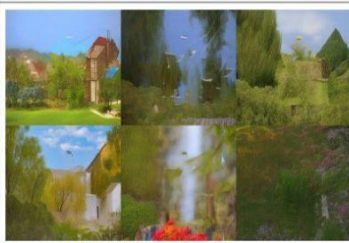
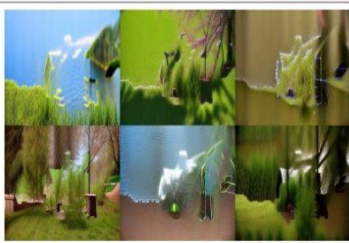
Insufficient prompt  
(w/o mentioning "house")  
"high-quality and detailed masterpiece"



Conflicting prompt  
"delicious cake"



Perfect prompt  
"a house, high-quality,  
extremely detailed, 4K, HQ"



# EVALUATION OF RESULTS

Method	Result Quality $\uparrow$	Condition Fidelity $\uparrow$
PITI [89](sketch)	$1.10 \pm 0.05$	$1.02 \pm 0.01$
Sketch-Guided [88] ( $\beta = 1.6$ )	$3.21 \pm 0.62$	$2.31 \pm 0.57$
Sketch-Guided [88] ( $\beta = 3.2$ )	$2.52 \pm 0.44$	$3.28 \pm 0.72$
ControlNet-lite	$3.93 \pm 0.59$	$4.09 \pm 0.46$
ControlNet	<b><math>4.22 \pm 0.43</math></b>	<b><math>4.28 \pm 0.45</math></b>

Table 1: Average User Ranking (AUR) of result quality and condition fidelity. We report the user preference ranking (1 to 5 indicates worst to best) of different methods.

Method	FID $\downarrow$	CLIP-score $\uparrow$	CLIP-aes. $\uparrow$
Stable Diffusion	6.09	0.26	6.32
VQGAN [19](seg.)*	26.28	0.17	5.14
LDM [72](seg.)*	25.35	0.18	5.15
PITI [89](seg.)	19.74	0.20	5.77
ControlNet-lite	17.92	0.26	6.30
ControlNet	15.27	0.26	6.31

Table 3: Evaluation for image generation conditioned by semantic segmentation. We report FID, CLIP text-image score, and CLIP aesthetic scores for our method and other baselines. We also report the performance of Stable Diffusion without segmentation conditions. Methods marked with “\*” are trained from scratch.



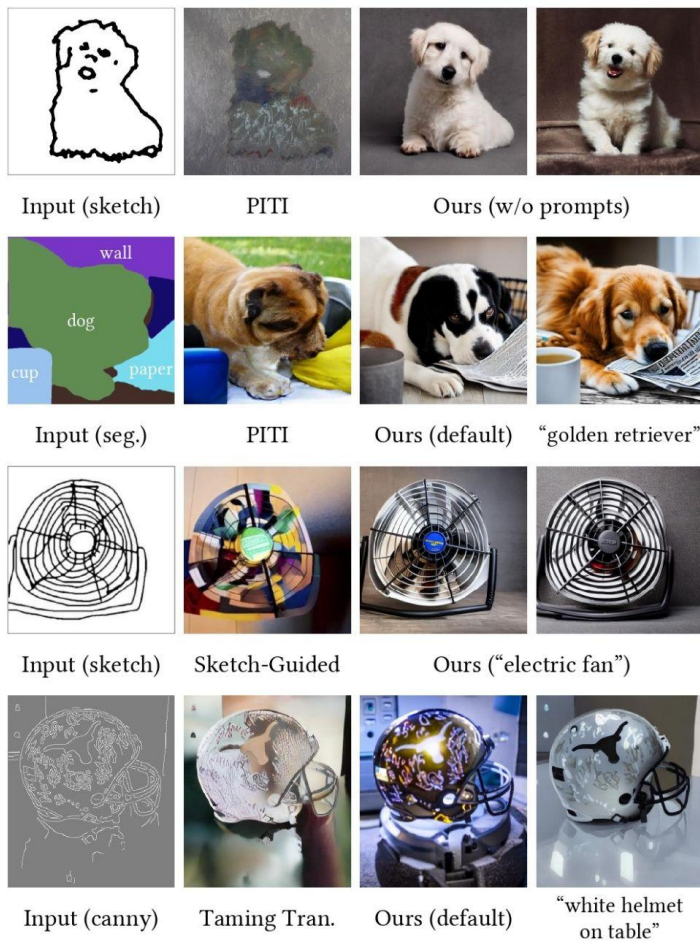


Figure 9: Comparison to previous methods. We present the qualitative comparisons to PITI [89], Sketch-Guided Diffusion [88], and Taming Transformers [19].

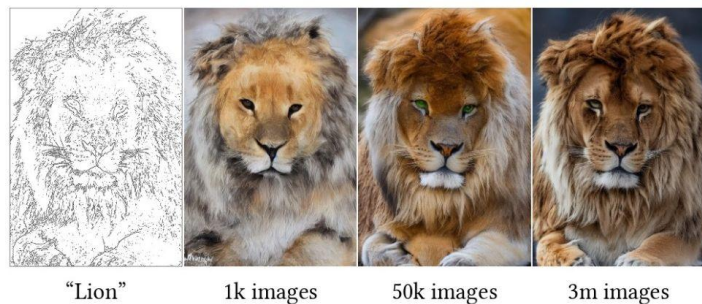


Figure 10: The influence of different training dataset sizes. See also the supplementary material for extended examples.



Figure 11: Interpreting contents. If the input is ambiguous and the user does not mention object contents in prompts, the results look like the model tries to interpret input shapes.



Figure 12: Transfer pretrained ControlNets to community models [16, 61] without training the neural networks again.

# USEFUL LINKS

paper: <https://arxiv.org/pdf/2302.05543>

paper2: <https://arxiv.org/pdf/2302.05543v1>

video: <https://youtu.be/WgrmCVa35ws?feature=shared>

hugging face: <https://huggingface.co/spaces/hysts/ControlNet>

github: <https://github.com/lllyasviel/ControlNet>

how to run in google collab:

<https://www.youtube.com/watch?v=Uq9N0nqUYqc>

Thanks for your attention!

