

# Suppress Content Shift: Better Diffusion Features via Off-the-Shelf Generation Techniques

Benyuan Meng, Qianqian Xu\*, Zitai Wang,  
Zhiyong Yang, Xiaochun Cao, Qingming Huang\*



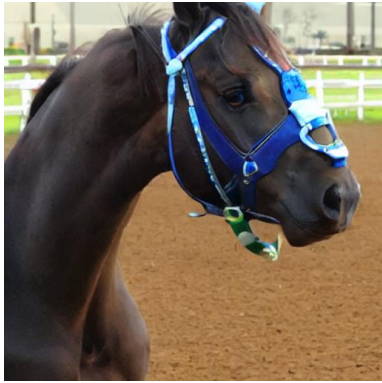
**Benyuan Meng**

2024.11

# Background

---

## Basic Generation



# Background

---

## Basic Generation



**Advanced Generation:  
Variation, Inpainting, Personalization...**

# Background

---

## Basic Generation

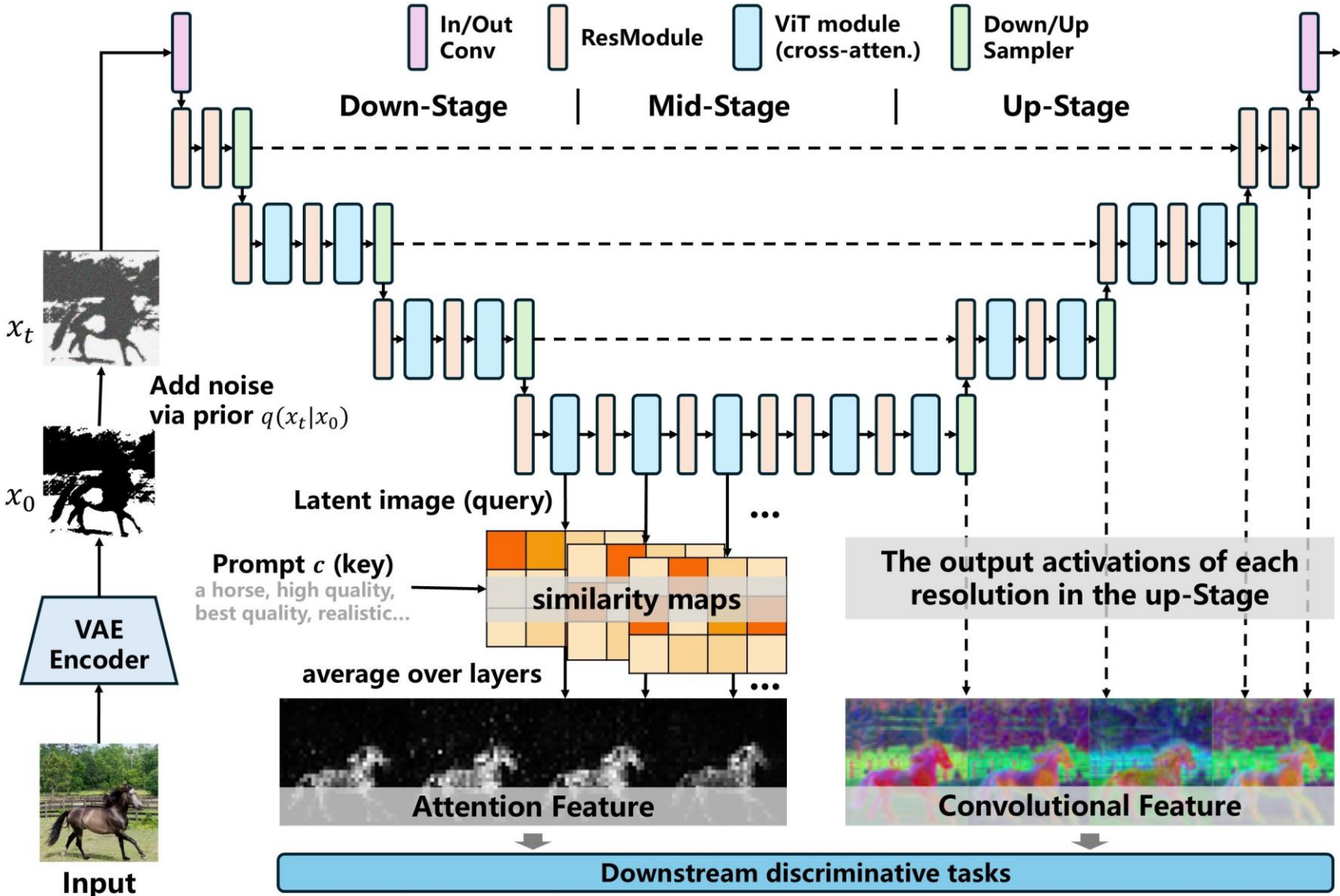


**Advanced Generation:  
Variation, Inpainting, Personalization...**

**Generation  $\rightarrow p(x, y) \rightarrow p(y|x) \rightarrow$  Discrimination**



# Background

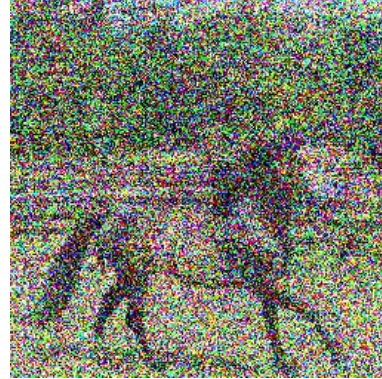


# Key Observation

---

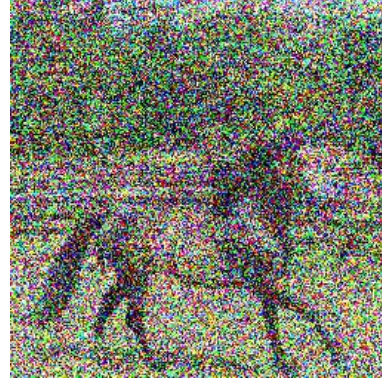


# Key Observation





# Key Observation





# Key Observation



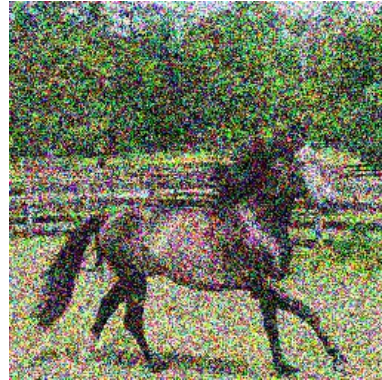
**Diffusion  
Model**

# Key Observation



×

Diffusion  
Model



✓

Diffusion  
Model



- Adding noises to input images is widely considered necessary to avoid singularity of clean inputs.



# Key Observation



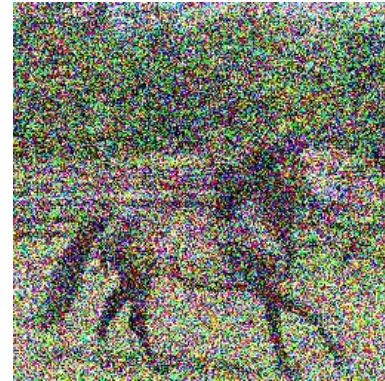
↓ ×

Diffusion  
Model



↓ ✓

Diffusion  
Model

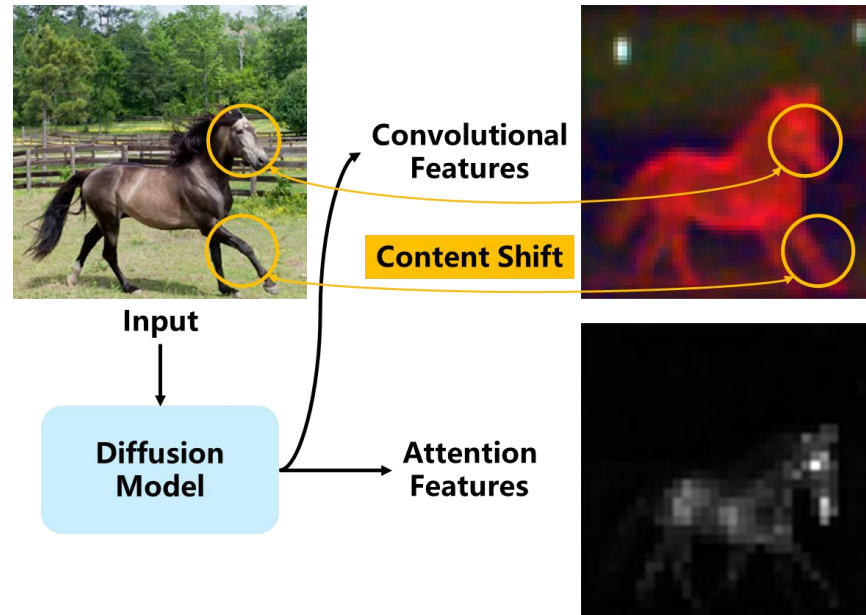








Information Loss?

- Adding noises to input images is widely considered necessary to avoid singularity of clean inputs.



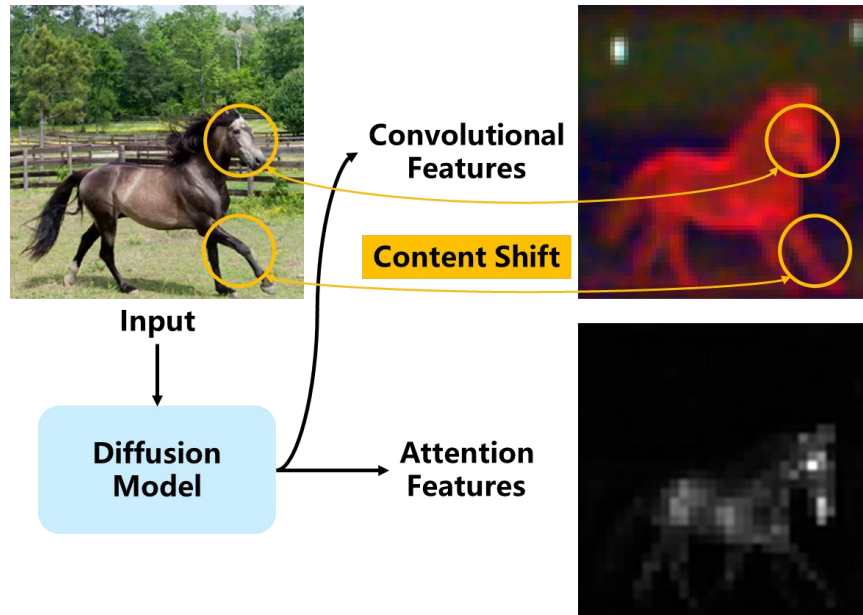
# Key Observation









Dataset	Low Quality	Neutral	High Quality
Horse-21	 58.90	 59.03	 59.33
CIFAR10	 91.67	 91.44	 91.21

- **Qualitative: perceivable content difference.**
- **Quantitative: non-negligible performance degradation.**

# Key Observation



Dataset	Low Quality	Neutral	High Quality
Horse-21	 58.90	 59.03	 59.33
CIFAR10	 91.67	 91.44	 91.21

- **Qualitative: perceivable content difference.**
- **Quantitative: non-negligible performance degradation.**
- **Content Shift**

# Key Observation

---



**Less Noises**  
**+ Less Content Shift**  
**- More Singularity**



**More Noises**  
**+ Less Singularity**  
**- More Content Shift**



# Key Observation



**Less Noises**  
**+ Less Content Shift**  
**- More Singularity**



**More Noises**  
**+ Less Singularity**  
**- More Content Shift**

**Additional  
Control  
Method**

**Better  
Feature**

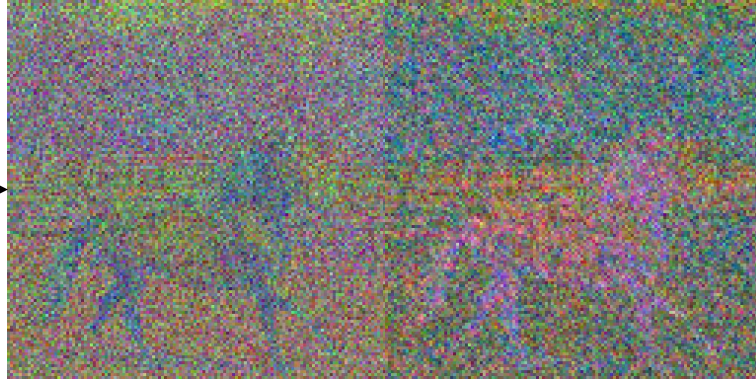
- **Not seeking balance.**
- **But a new way to suppress content shift while maintaining noise strength.**

# Understanding Content Shift

---



**Original Image**



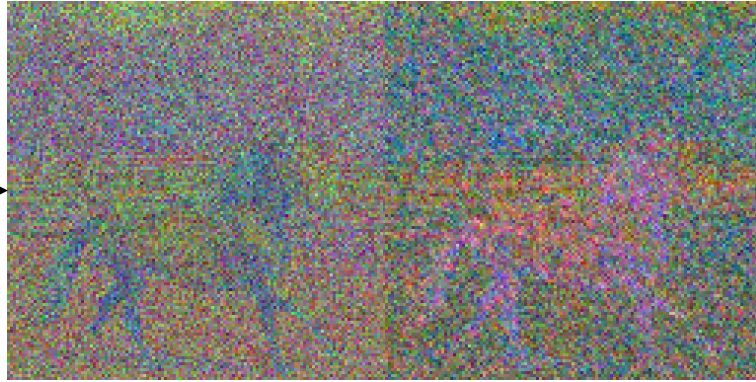
**Noisy Activations  
in Early Sections**

# Understanding Content Shift

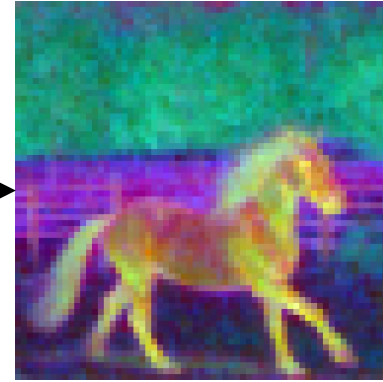
---



**Original Image**



**Noisy Activations  
in Early Sections**



**Reconstructed  
Clean Representations**

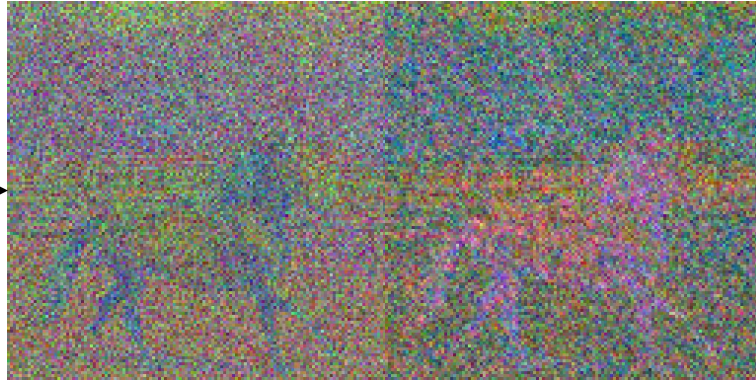


# Understanding Content Shift

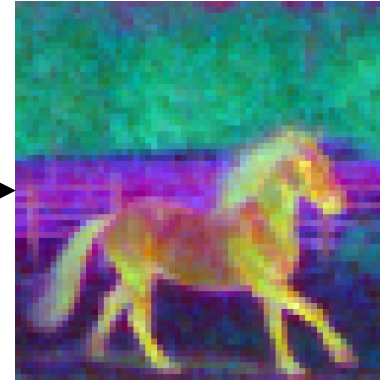
---



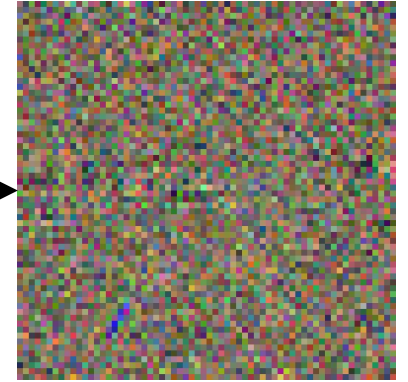
**Original Image**



**Noisy Activations  
in Early Sections**



**Reconstructed  
Clean Representations**



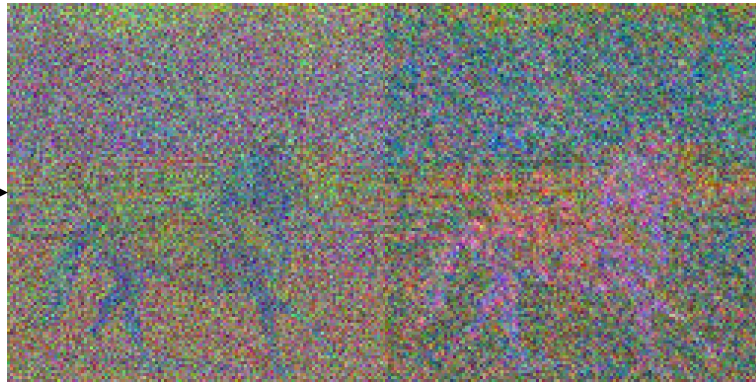
**Output Noises**

# Understanding Content Shift

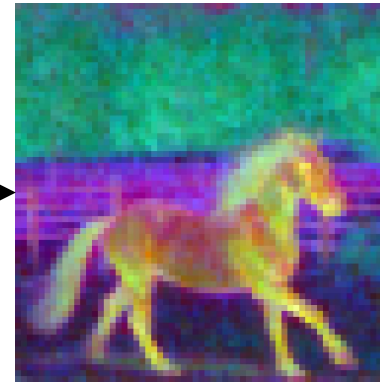
---



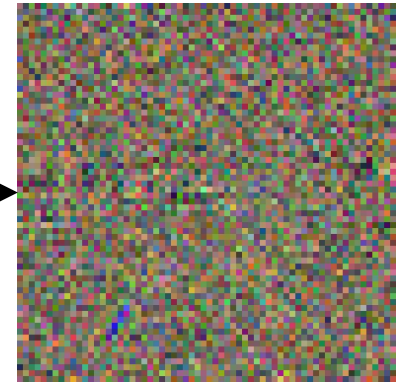
Original Image



Noisy Activations  
in Early Sections



Reconstructed  
Clean Representations

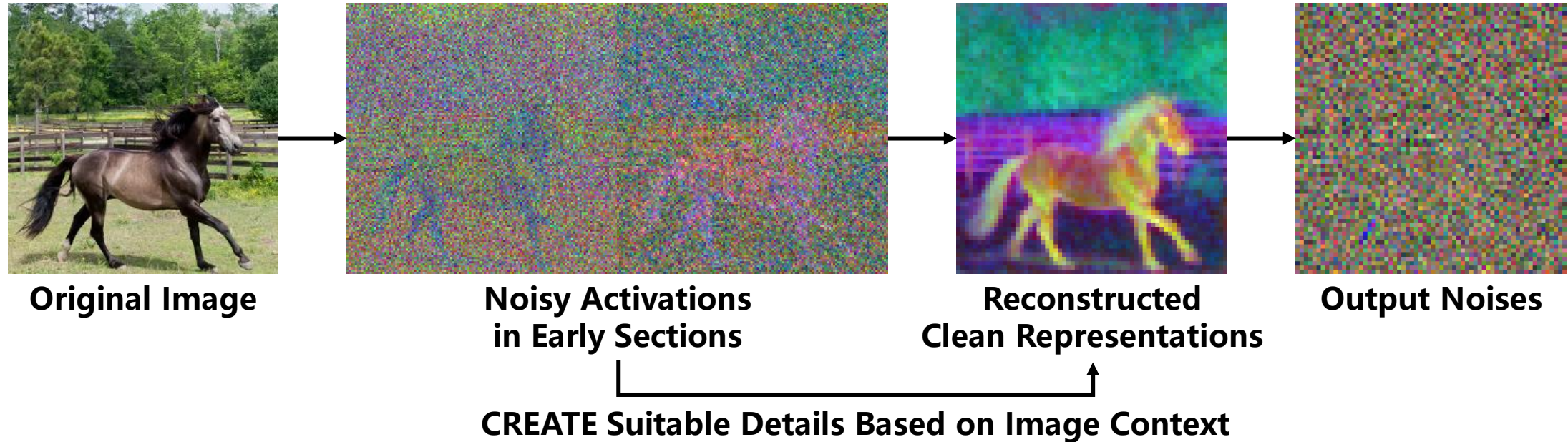


Output Noises

- Diffusion models have learnt to reconstruct clean representations from noisy inputs.



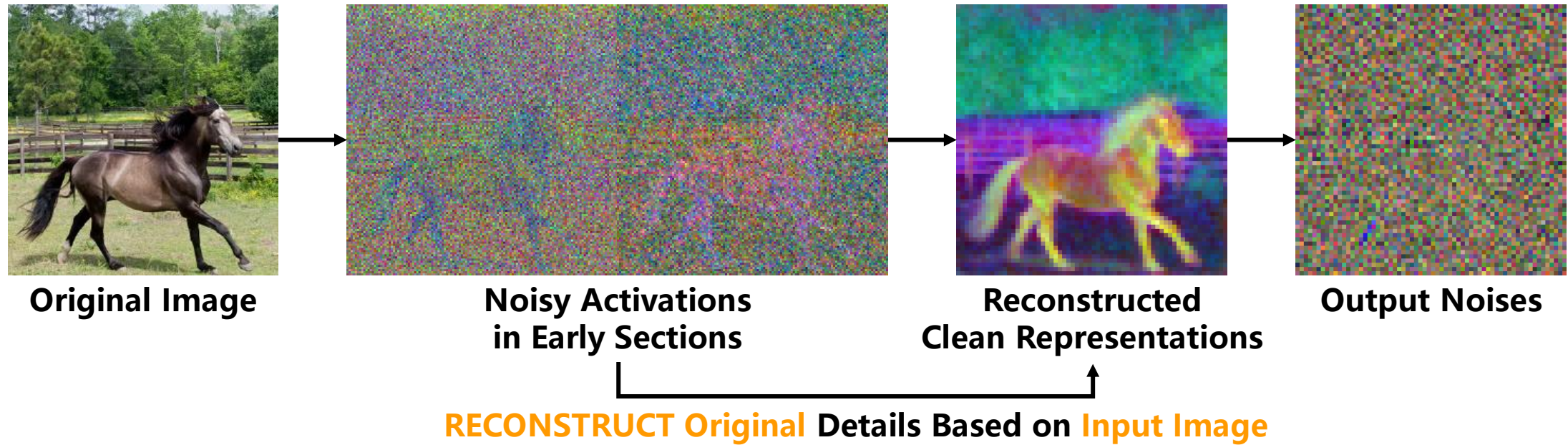
# Understanding Content Shift



- Diffusion models have learnt to reconstruct clean representations from noisy inputs.
- Reconstructed details are created based on context, not guaranteed to stay true to the original image.

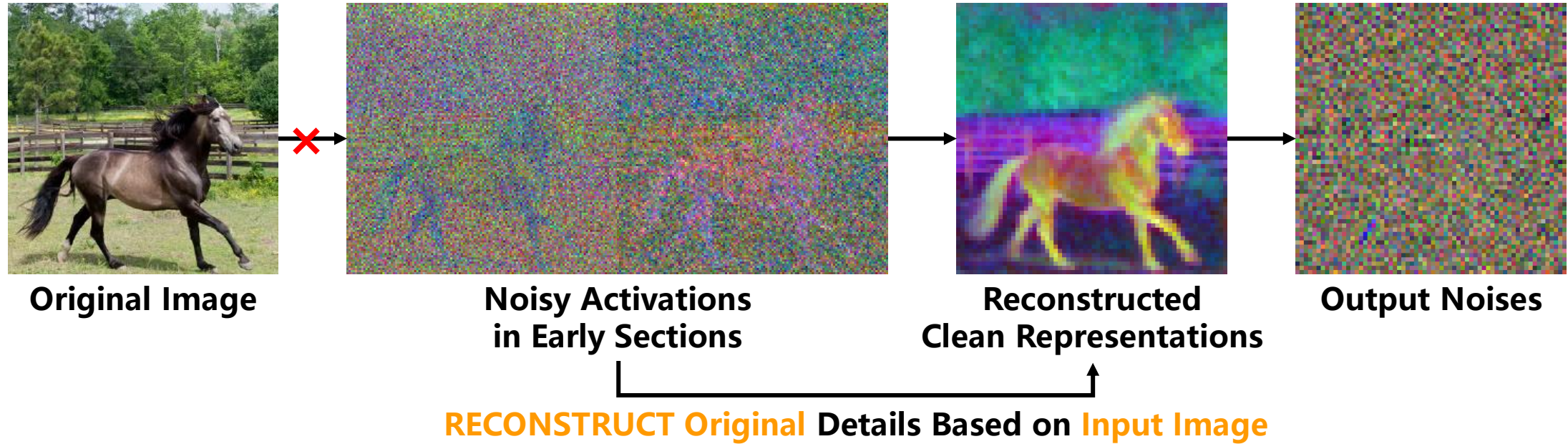


# Suppressing Content Shift



- Steer reconstruction?

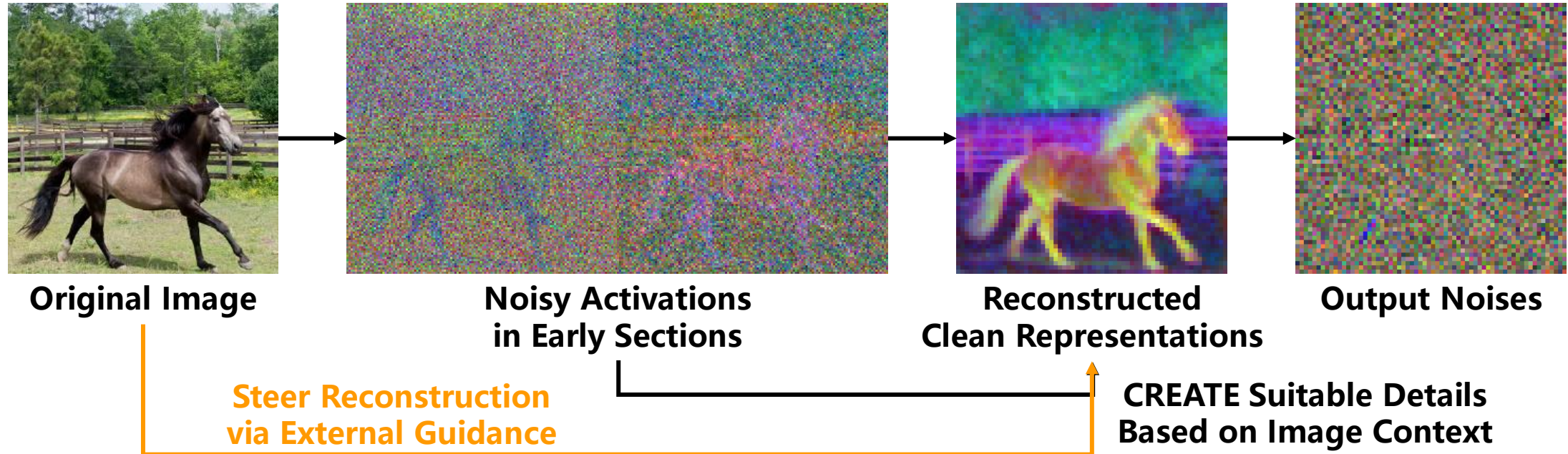
# Suppressing Content Shift



- **Steer reconstruction?**



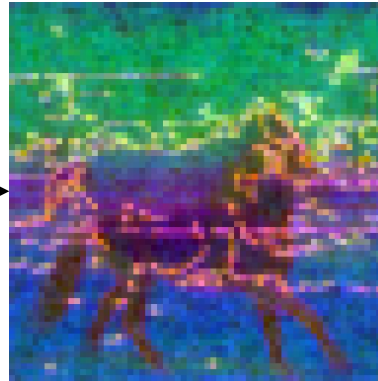
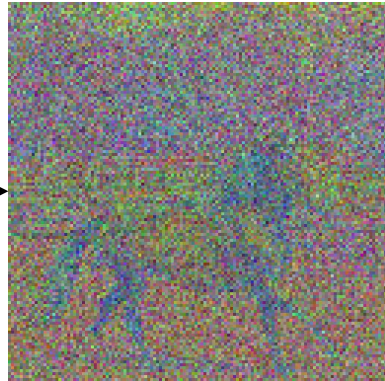
# Suppressing Content Shift



- **Steer reconstruction?**
- **Bypass the normal route to directly inject guidance information.**



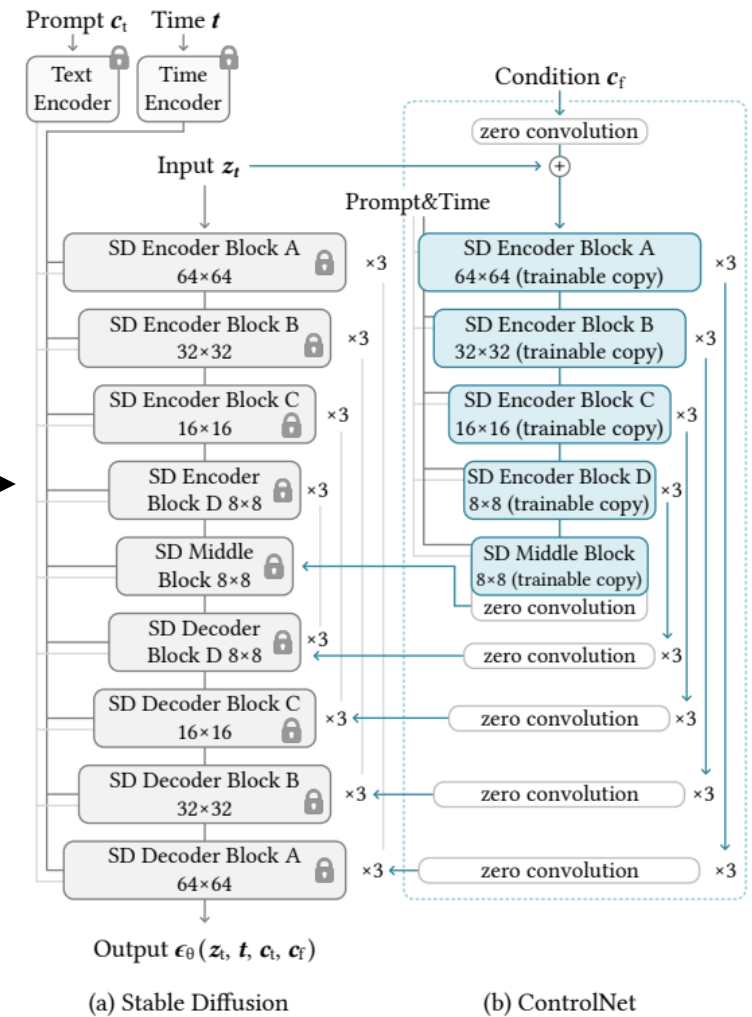
# Suppressing Content Shift



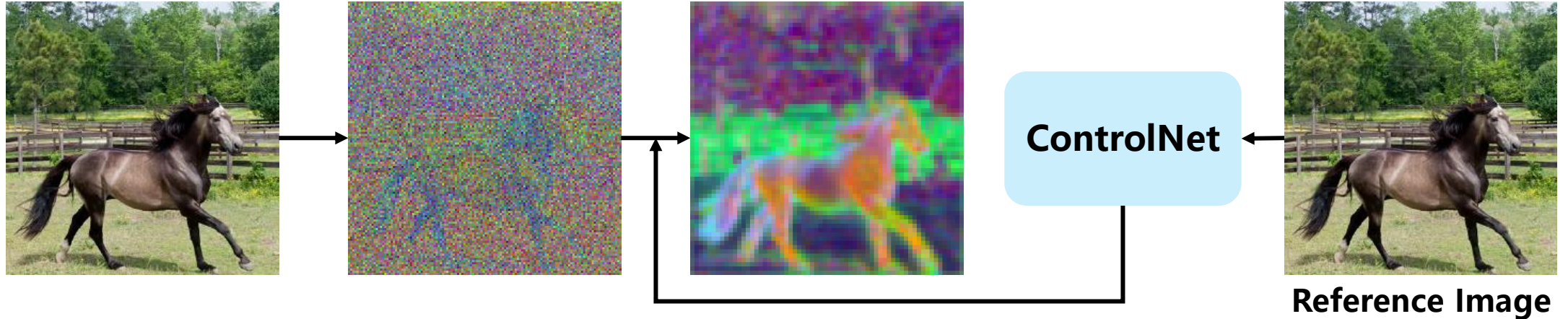
**ControlNet**

Reference Image

Taken from Zhang, Lvmin, Anyi Rao, and Maneesh Agrawala. "Adding conditional control to text-to-image diffusion models." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023.

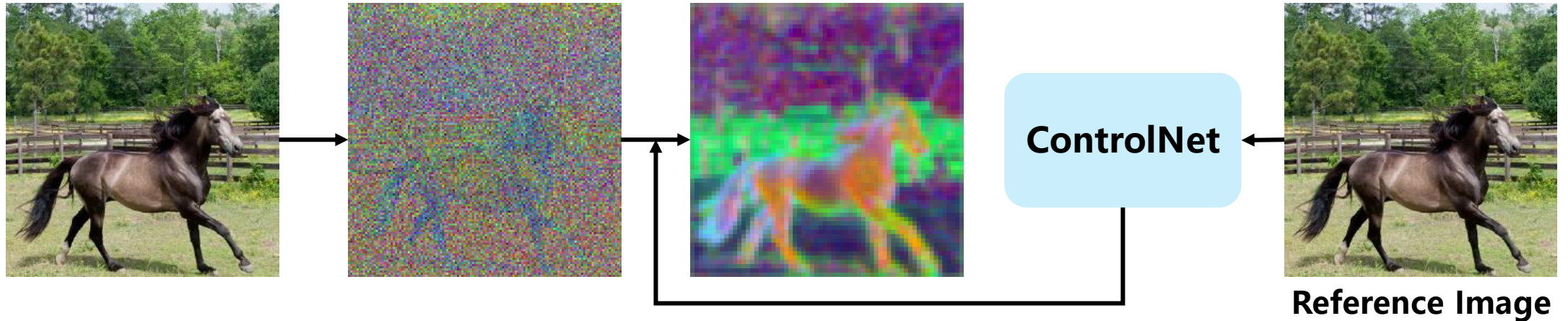


# Suppressing Content Shift



- **Utilize generation control techniques to steer reconstruction to stay close to the original image.**
- **Named GATE**
  - **(GenerAtion Techniques Enhanced diffusion feature)**

# Suppressing Content Shift



- **Utilize generation control techniques to steer reconstruction to stay close to the original image.**
  - Fine-Grained Prompts
  - ControlNet
  - LoRA
- **Feature amalgamation to harness enhanced diversity.**



# Results

Table 1: The results of semantic correspondence (left, PCK@0.1) and label-scarce semantic segmentation (right, mIoU $\uparrow$ ). **Red** for the best result and **blue** for the runner-up.

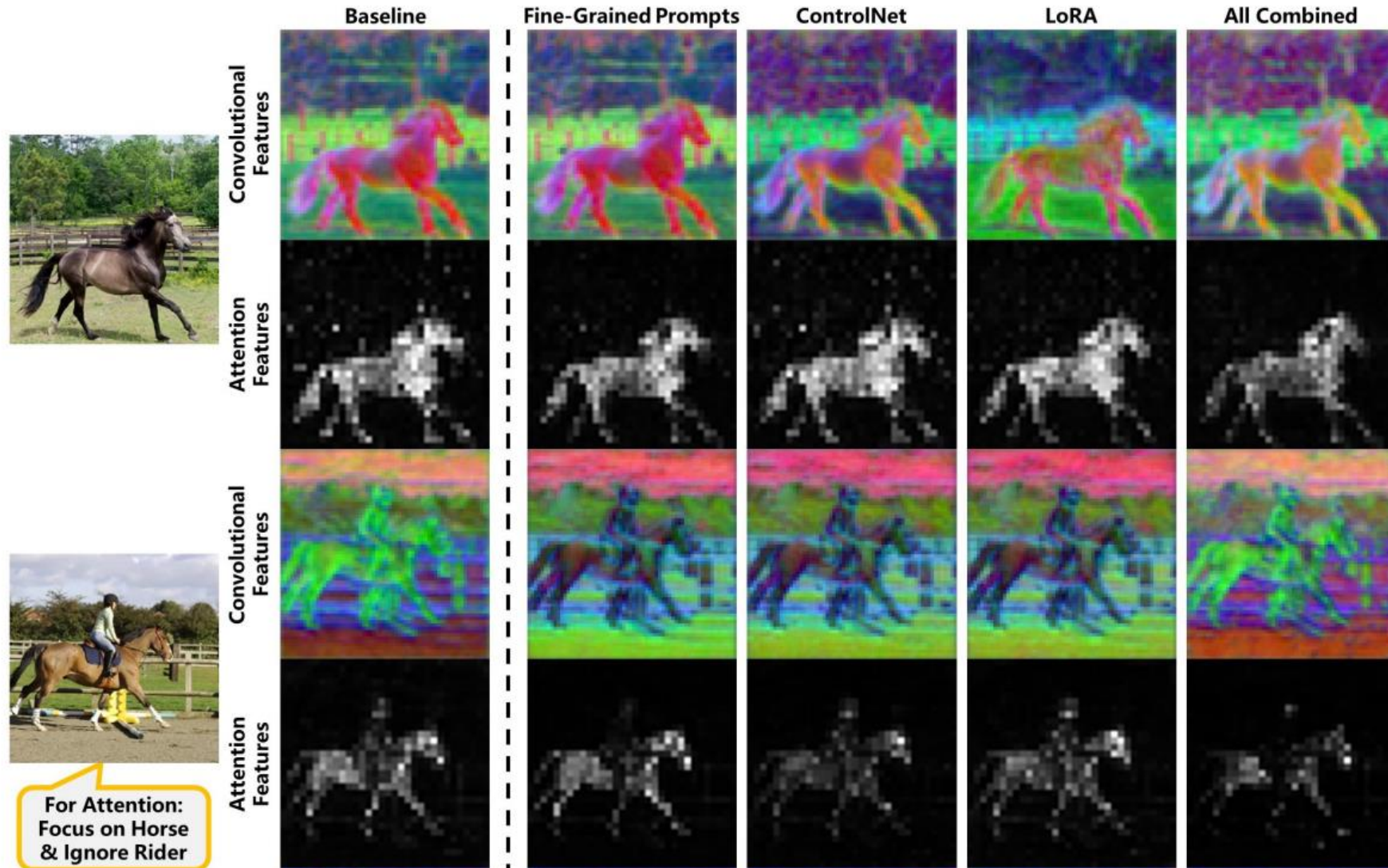
Method		PCK <sub>img</sub> $\uparrow$	PCK <sub>bbox</sub> $\uparrow$	Method	Bedroom-28	Horse-21
Non-DF	DINO	51.68	41.04	ALAE	20.0 $\pm$ 1.0	–
	DHPF	55.28	42.63	GAN Inversion	13.9 $\pm$ 0.6	17.7 $\pm$ 0.4
DF	DIFT	-	52.90	GAN Encoder	22.4 $\pm$ 1.6	26.7 $\pm$ 0.7
	DHF	72.56	64.61	SwAV	41.0 $\pm$ 2.3	51.7 $\pm$ 0.5
Baseline	nn	61.15	51.66	SwAVw2	42.4 $\pm$ 1.7	54.0 $\pm$ 0.9
	conv	<b>73.96</b>	<b>65.74</b>	MAE	45.0 $\pm$ 2.0	63.4 $\pm$ 1.4
<b>GATE</b>	nn	64.47	55.72	DatasetGAN	31.3 $\pm$ 2.7	45.4 $\pm$ 1.4
	conv	<b>76.60</b>	<b>69.10</b>	DatasetDDPM	47.9 $\pm$ 2.9	60.8 $\pm$ 1.0
				DDPM	<b>49.4 <math>\pm</math> 1.9</b>	<b>65.0 <math>\pm</math> 0.8</b>
				<b>GATE</b>	<b>53.1 <math>\pm</math> 2.7</b>	<b>67.2 <math>\pm</math> 1.1</b>

# Results

Table 2: Results on the two standard semantic segmentation datasets, ADE20K and CityScapes. **Red** for the best result and **blue** for the runner-up.

Category	Method	ADE20K			CityScapes		
		mIoU $\uparrow$	aAcc $\uparrow$	mAcc $\uparrow$	mIoU $\uparrow$	aAcc $\uparrow$	mAcc $\uparrow$
SOTA	MaskCLIP	23.70	-	-	-	-	-
	ODISE	29.90	-	-	-	-	-
	VPD	<b>37.63</b>	<b>79.16</b>	<b>50.08</b>	<b>55.06</b>	<b>90.14</b>	<b>68.96</b>
Ours	GATE	<b>40.51</b>	<b>79.68</b>	<b>54.90</b>	<b>64.20</b>	<b>92.83</b>	<b>76.98</b>

# Results





---

**Thanks for your listening!**

