# NEP: Autoregressive Image Editing via Next Editing Token Prediction

Huimin Wu<sup>1</sup>, Xiaojian Ma<sup>1</sup>, Haozhe Zhao<sup>2</sup>, Yanpeng Zhao<sup>1</sup>, Qing Li<sup>1</sup>

<sup>1</sup>BIGAI <sup>2</sup>Peking University

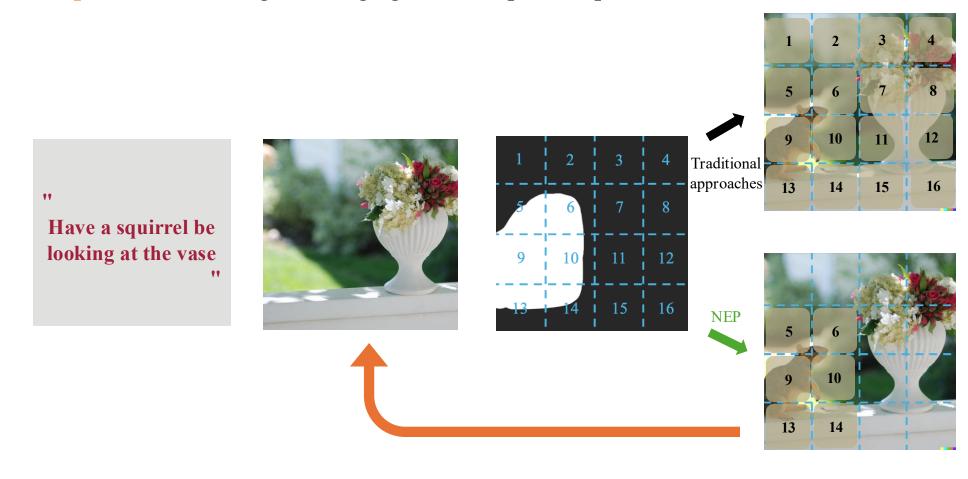
#### Text-driven Image Editing



#### We propose NEP, a targeted and efficient autoregressive image editing method.

It allows regeneration solely at edited regions.

It makes **self-improvement** during the image generation process possible.

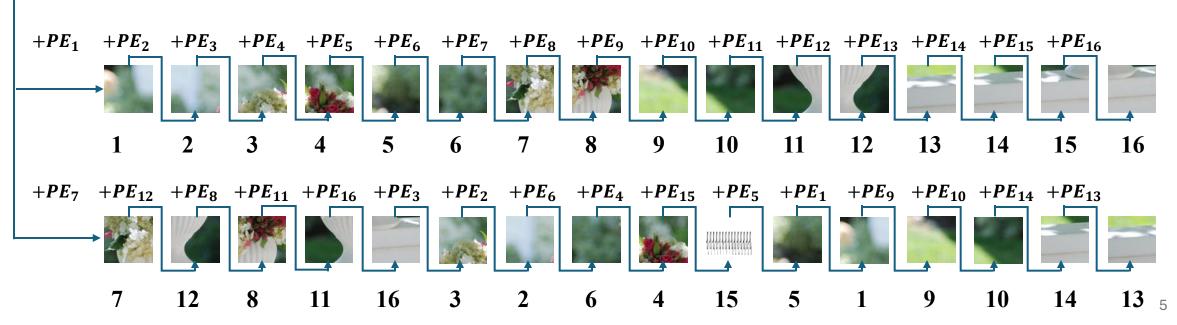


- 1) NEP Pre-training
- 2) NEP fine-tuning

- **5**1
- ) NEP Pre-training
- 2) NEP fine-tuning

A flower vase is sitting on a porch stand





- •
- NEP Pre-training
- 2) NEP fine-tuning



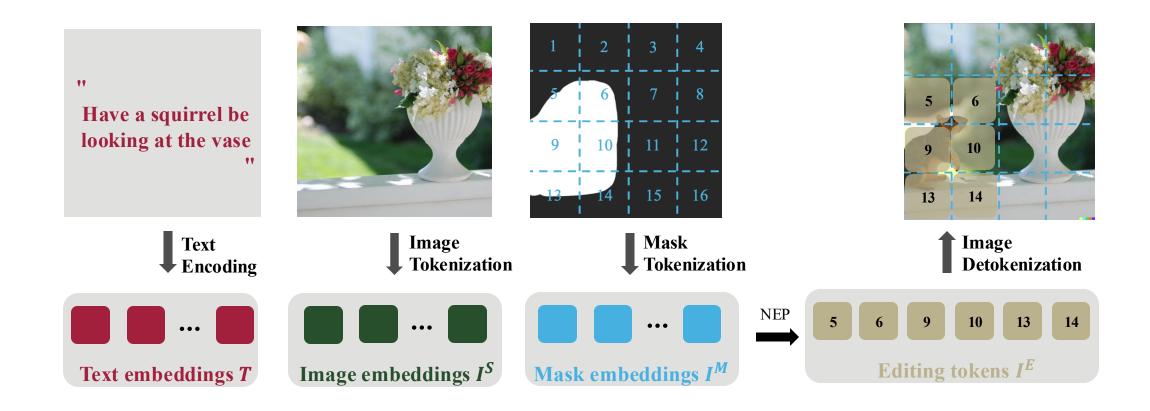
"Pancake with brown topping and ice cream."

"Club sandwich with fries and mustard."

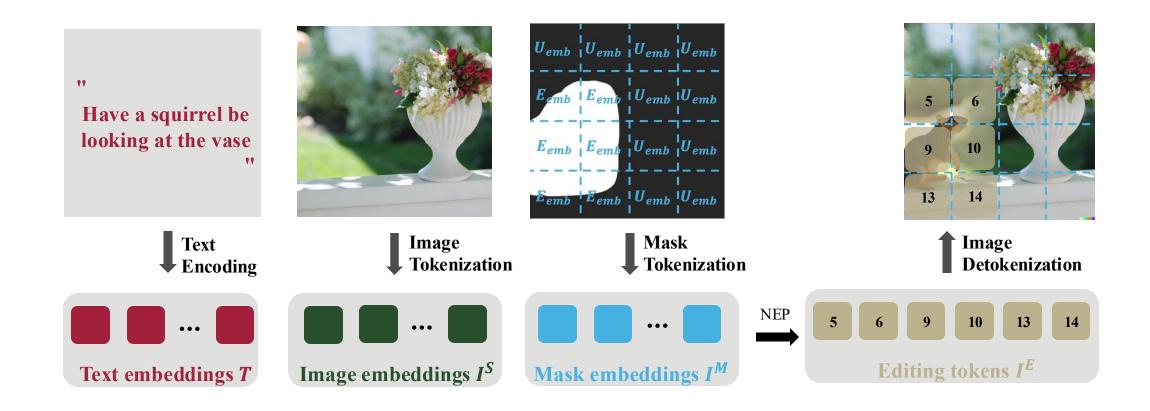
"A black refrigerator with its cabinet door closed in a kitchen."

"Red surfboard on grass in backyard."

- 1) NEP Pre-training
- 2) NEP fine-tuning



- 1) NEP Pre-training
- 2) NEP fine-tuning



Settings	Methods	L1↓	L2↓	CLIP-I↑	DINO↑
	Global Description-guided				
Single-turn	SD-SDEdit Null Text Inversion	0.1014 0.0749	0.0278 0.0197	0.8526 0.8827	0.7726 0.8206
	GLIDE Blended Diffusion	3.4973 3.5631	115.8347 119.2813	0.9487 0.9291	0.9206 0.8644
	Instruction-guided				
	HIVE InstructPix2Pix (IP2P)	0.1092 0.1141	0.0380 0.0371	0.8519 0.8512	0.7500 0.7437
	IP2P w/ MagicBrush UltraEdit	0.0625 0.0575	0.0203 $0.0172$	$\frac{0}{0.9332}$	$\frac{0.8987}{0.8982}$
	FireEdit AnySD	0.0701 0.111 <u>4</u>	0.0238 0.0439	0.9131 0 <u>.8</u> 676	0.8619 0.76 <u>8</u> 0
	Ours	0.0547	0.0163	0.9350	0.9044
Multi-turn	Global Description-guided				
	SD-SDEdit Null Text Inversion	0.1616 0.1057	0.0602 0.0335	0.7933 0.8468	0.6212 0.7529
	GLIDE Blended Diffusion	11.7487 14.5439	1079.5997 1510.2271	0.9094 0.8782	0.8494 0.7690
	Instruction-guided				
	HIVE InstructPix2Pix (IP2P)	0.1521 0.1345	0.0557 0.0460	0.8004 0.8304	0.6463 0.7018
	IP2P w/ MagicBrush UltraEdit	0.0964 0.0745	0.0353 <b>0.0236</b>	0.8924 0.9045	0.8273 0.8505
	FireEdit AnySD	0.0911 0.0748	0.0326 0.0273	0.8819 <b>0.9152</b>	0.8010 <b>0.8623</b>
	Ours	0.0707	<u>0.0269</u>	<u>0.9107</u>	0.8493

Input Mask Ground truth IP2P MagicBrush UltraEdit AnySD NEP (ours)



"Could it be a glass of wine on the table?"



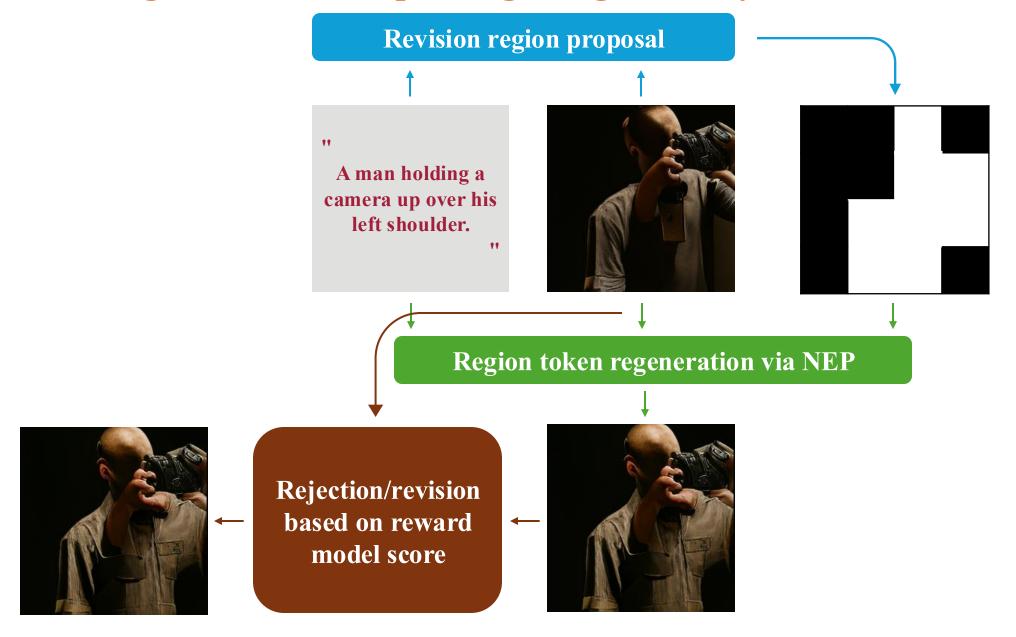
"let there be sports on TV"



"let the woman wear a designer gown"

## Test-time Scaling of NEP for Improving Image Quality

#### **Test-time Scaling of NEP for Improving Image Quality**



#### **Test-time Scaling of NEP for Improving Image Quality**



"A man holding a camera up over his left shoulder."



"A passenger plane that is parked on the runway."



"A boat floating on a lake through a brick bridge."



"A kitchen filled with lots of counter top space."

# THANK YOU!