

# Visual Instruction Inversion: Image Editing via Visual Prompting — *Supplementary Material* —

Anonymous Author(s)

Affiliation

Address

email

1 This document provides additional information complementing the main paper. First, we compare  
2 to Textual Inversion in Sec. A. Then, in Sec. B, we provide additional qualitative comparisons to  
3 Imagic [3] and Null-text Inversion [5]. Finally, in Sec. C, we provide the implementation details of  
4 our method, along with results obtained using a variant of our approach - instruction concatenation.  
5 This variant allows users to add extra information into the learned instruction.

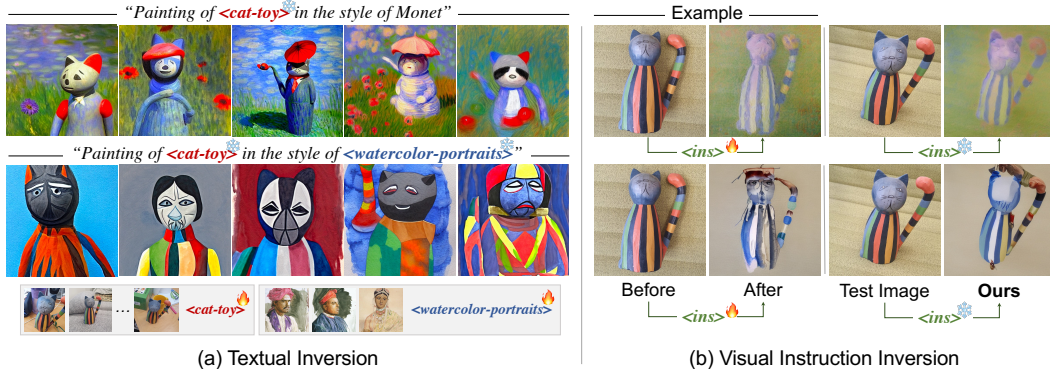


Figure 1: **Ours vs. Textual Inversion.** (a) Textual Inversion inverts a visual concept or object (e.g., a particular <cat-toy>) into a word embedding. This optimized word embedding can then be combined with a textual description to generate novel scenes. (b) Our Visual Instruction Inversion learns the transformation <ins>: "before" → "after" in a given before-and-after image pair. This learned instruction can then be applied to new test images to perform the same edit.

## 6 A Textual Inversion vs. Visual Instruction Inversion

7 Textual Inversion [2, 6] is a method to invert a visual concept into a corresponding representation in  
8 the language space. In particular, given (i) a text-to-image pre-trained model and (ii) some images  
9 describing a visual concept (e.g., a particular kind of toy; Figure 1 bottom row), Textual Inversion  
10 learns new "words" in the embedding space of the text-to-image model to represent those visual  
11 concepts. Once these "words" are learned for that concept, they can be plugged into arbitrary  
12 textual descriptions, just like other English words, which can then be used to create the target visual  
13 concept in different contexts. Instead of learning the representation for an isolated visual concept, our  
14 approach (Visual Instruction Inversion), learns the *transformation* from a before-and-after image pair.  
15 This learned transformation is then applied to a test image to achieve similar edit "before" → "after".

16 **Applicability of Textual Inversion for image editing.** Given these differences with our proposed  
17 method, we now try to see if Textual Inversion can be used for image editing. Textual Inversion can

18 generate a “painting of a <cat-toy> in the style of Monet” by using the learned word <cat-toy>  
19 from example images (Figure 1a, Row 1). However, the synthesized images often only capture the  
20 essence of the objects, and disregard the details of the input images. As a result, textual inversion is  
21 suitable for novel scene composition, but is not effective for image editing.

22 On the other hand, our Visual Instruction Inversion does not learn novel token representations for  
23 objects or concepts. Instead, we learn the edit instruction from before-and-after pairs, which can be  
24 applied to any test image to obtain corresponding edits. This allows us to achieve fine-grained control  
25 over the resulting images. For example, by providing a photo of <cat-toy>, one before and one in  
26 a specific impressionist style, we learn the transformation from before to impressionist, denoted as  
27 <ins>. Once learned, this instruction can be applied to new <cat-toy> images to achieve the same  
28 impressionist painting style, without losing the fine details of the test image (Figure 1b, Row 1).

29 One might suggest an alternative approach to image editing using Textual Inversion, which involves  
30 learning two tokens: one for the object and another for the style (e.g., “Painting of <cat-toy> in the  
31 style of <watercolor-portraits>”). Figure 1a (Row 2) shows the results of this approach. As  
32 can be seen, Textual Inversion still often introduces significant changes that deviate from the original  
33 input image. Thus, Textual Inversion is not suitable for accurate image editing.

## 34 B Additional Qualitative Comparisons

35 In the attached HTML file (index.html), we provide additional comprehensive qualitative compar-  
36 isons for our method versus InstructPix2Pix [1] and SDEdit [4].

37 We also present qualitative comparisons with other state-of-the-art text-conditioned image editing  
38 methods, Imagic [3] and Null-text Inversion [5] (Figure 2). These methods can generate outputs based  
39 on given text prompts, such as “A watercolor painting of a cat” (Row 3). However, the outputs often  
40 do not match the given reference. The text prompts can also be ambiguous and result in unsatisfactory  
41 outputs, as illustrated by the case of “A character in a Pixar movie” (Row 1).

42 Another challenge is the inconsistency of text-conditioned models, where the same text prompt  
43 can produce different outputs for different test images. For example, the text prompt “A frozen  
44 waterfall” (Row 6) generates different water colors (blue vs. white) when applied to different test  
45 images (Before-and-after pair is from [5]). Our method is more consistent in this case, as the learned  
46 instruction might have learned the water color.

## 47 C Implementation Details

### 48 C.1 Optimization settings

49 We use the pretrained clip-vit-large-patch14 as the CLIP Encoder in our approach. For  
50 instruction initialization [7], we set the caption length for after image to 10 tokens. However, this  
51 specific caption length does not affect the optimization algorithm. We can optimize initialization  
52 instructions of varying lengths (up to 77 tokens). It takes roughly 7 minutes to optimize for one edit,  
53 and 4 seconds to apply the learned instruction to new images.

54 Specifically, during the optimization process, we freeze the tokens representing the start of text  
55 (<|startoftext|>), end of text (<|endoftext|>), and all padding tokens after end of text  
56 (<|endoftext|>). We only update the tokens inside the text prompt, called <ins> (between  
57 <|startoftext|> and <|endoftext|>) (Figure 3a).

### 58 C.2 Instruction Concatenation

59 We only optimize a fixed number of tokens, so we have the flexibility to concatenate additional  
60 information to the learned instruction during inference (Figure 3b). This allows us to achieve more  
61 fine-grained control over the resulting images. Figure 4 shows qualitative results of this approach.  
62 Users can input extra information to combine or guide the learned instruction according to their  
63 preferences.

64 In the first example, we transform “cat” → “watercolor cat”. We demonstrate how concatenating extra  
65 information to the learned instruction enables both image editing (changing to a watercolor style) and

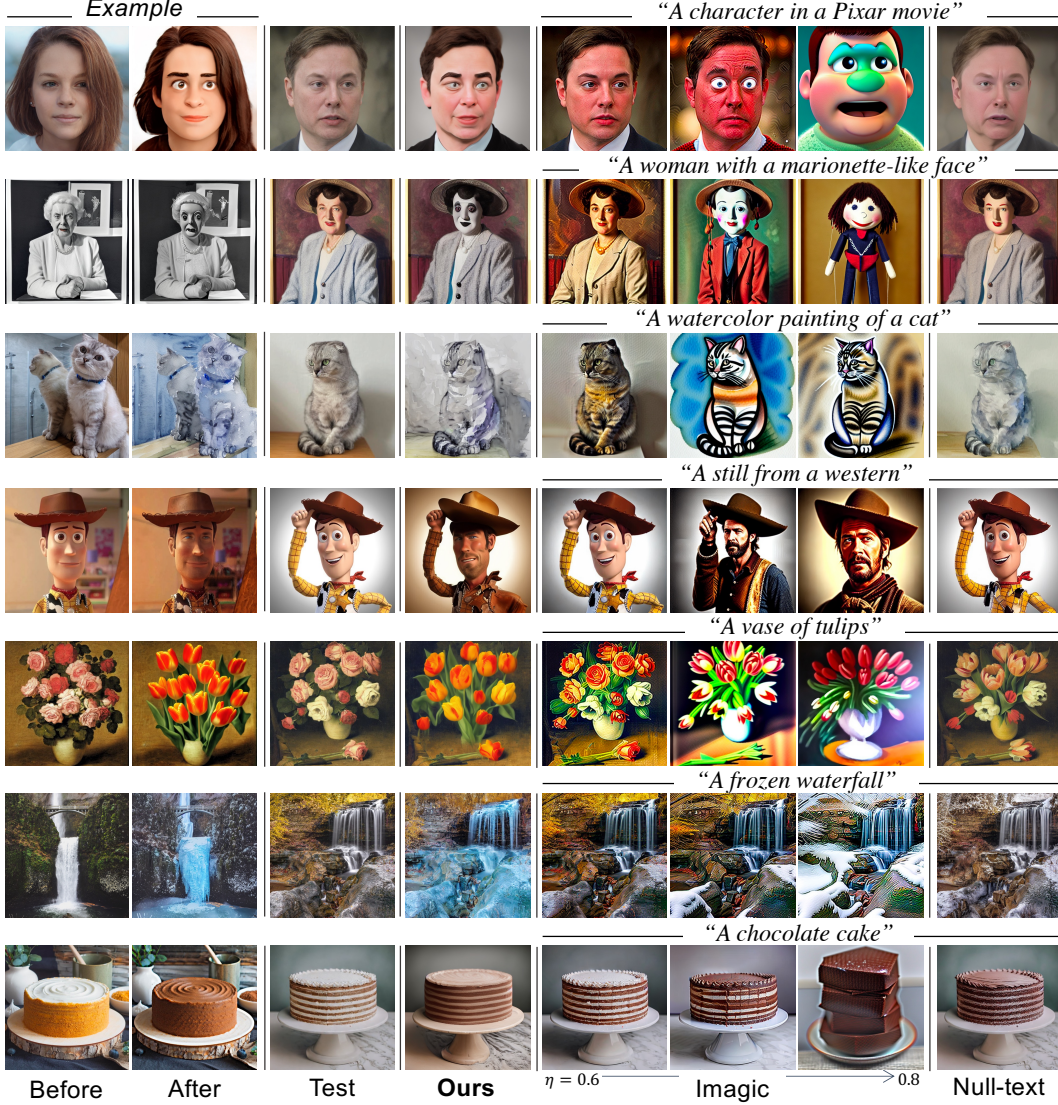


Figure 2: **Additional qualitative comparisons** to Imagic [3] and Null-text Inversion [5]. Imagic and Null-text Inversion fail to match the reference image as they perform edits based on ambiguous text prompts (Row 1-4); or exhibit inconsistency in producing outputs for the same prompt across test images (Row 6). In contrast, our method produces visually closer edited images to the before-and-after pair while demonstrating improved consistency by using the learned instructions.

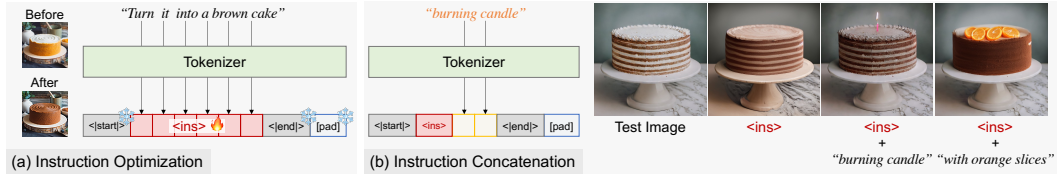


Figure 3: **Implementation details.** (a) Instruction Optimization: We only optimize a part of the instruction embedding, called `<ins>`. (b) Instruction Concatenation: During test time, we can add extra information into the learned instruction `<ins>` to further guide the edit.

66 domain translation (e.g., “cat” → “tiger”). The painting style is consistent with the before-and-after  
 67 images, while the domain translation corresponds to the additional information provided by the user.



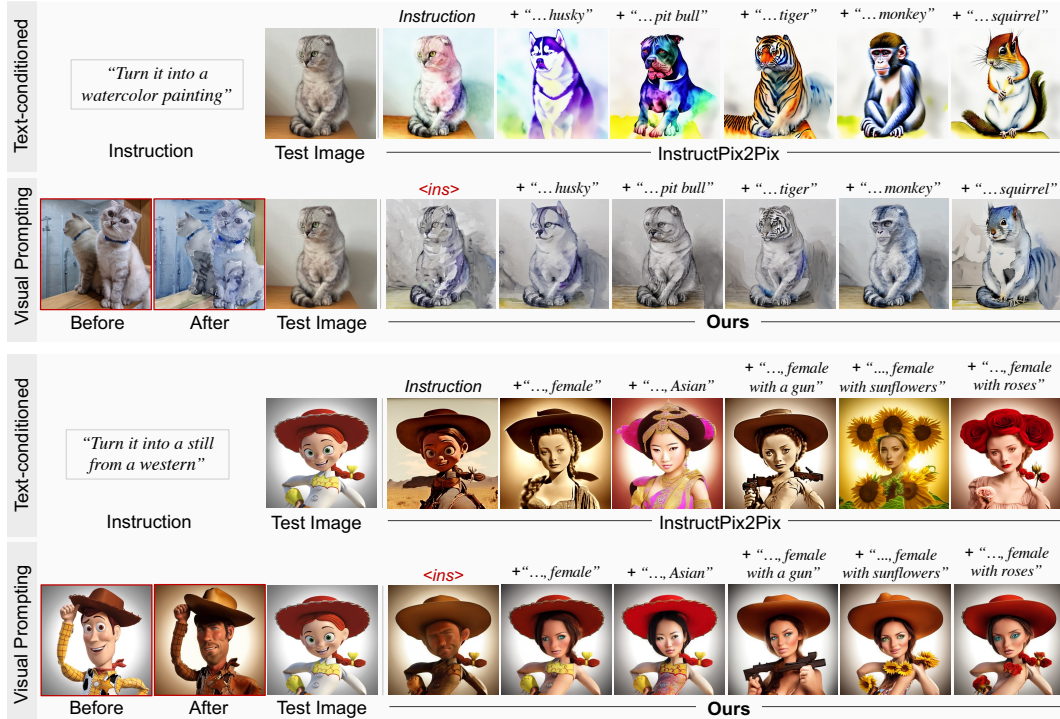


Figure 4: **Instruction concatenation.** We can concatenate extra information into the learned instruction `<ins>` to navigate the edit. (Zoom in for details).

68 Applying InstructPix2Pix [1] often does not yield satisfactory results, as the painting style differs  
69 from the reference image.

70 In the second example, we further illustrate how concatenating extra information can help guide  
71 the learned instruction. We initially learned a male character transformation from a cartoon image.  
72 However, applying this learned instruction to a new female test image results in a male bias in the  
73 output. By adding extra information, we can navigate the edit to overcome this bias. We can also  
74 add extra details such as sunflowers or guns to adjust the learned instruction accordingly. These  
75 modifications ensure consistent outputs that are aligned with the original learned instruction. Applying  
76 InstructPix2Pix, again, tends to produce more varied outputs that are inconsistent with the user’s  
77 intention.

## 78 Photo Attribution

- 79 • Elsa (Human): reddit.com/r/Frozen
- 80 • Disney characters: princess.disney.com
- 81 • Toy Story characters: toystory.disney.com
- 82 • Toonify faces: toonify.photos
- 83 • Girl with a Pearl Earring: wikipedia/girl-with-a-pearl-earring
- 84 • Mona Lisa: wikipedia/mona-lisa
- 85 • The Princesse de Broglie: wikipedia/Princesse-de-Broglie
- 86 • Self-portrait in a Straw Hat: wikipedia/self-portrait-in-a-straw-hat
- 87 • <cat-toy> and <watercolor-portraits> concept: huggingface.co/sd-concepts-library
- 88 • Gnoch cat, waterfall, and cake images are from Imagic [3] and Null-text Inversion [5].



## 89 References

- 90 [1] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing  
91 instructions. In *arXiv*, 2023.
- 92 [2] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-  
93 Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *arXiv*,  
94 2022.
- 95 [3] Bahjat Kavar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal  
96 Irani. Imagic: Text-based real image editing with diffusion models. In *Conference on Computer Vision and*  
97 *Pattern Recognition 2023*, 2023.
- 98 [4] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit:  
99 Guided image synthesis and editing with stochastic differential equations. In *arXiv*, 2022.
- 100 [5] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing  
101 real images using guided diffusion models. In *arXiv*, 2022.
- 102 [6] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dream-  
103 booth: Fine tuning text-to-image diffusion models for subject-driven generation. In *arXiv*, 2022.
- 104 [7] Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Hard  
105 prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery. In *arXiv*, 2023.