

TEXT-BASED VECTOR SKETCH EDITING WITH IMAGE EDITING DIFFUSION PRIOR –SUPPLEMENTAL DOCUMENT–

Haoran Mo, Xusheng Lin, Chengying Gao* and Ruomei Wang

Sun Yat-sen University

{mohaor,linxsh8}@mail2.sysu.edu.cn, {mcsgcy,isswrm}@mail.sysu.edu.cn

1. IMPLEMENTATION DETAILS

1.1. Iterative Editing

We show the implementation details of iterative editing in Algorithm 1, following some formulations in Prompt-to-Prompt [1].

Algorithm 1 Algorithm for iterative editing.

Input: All prompts $\{P^1, P^2, \dots, P^K\}$, a random seed ξ , and words $\{(w^1, w^2)^2, (w^2, w^3)^3, \dots, (w^{K-1}, w^K)^K\}$ specifying the editing region for local editing.

Output: All images x^1, x^2, \dots, x^K .

```
1:  $z_T \sim N(0, I)$ , a unit Gaussian random variable with  $\xi$ ;  
2:  $(z_T^1, z_T^2, \dots, z_T^K) \leftarrow z_T$ ;  
3: for  $t = T, T - 1, \dots, 1$  do  
4:    $z_{t-1}^1, M_t^1 \leftarrow DM(z_t^1, P^1, t, \xi)$ ;  
5:   for  $k = 2, \dots, K$  do  
6:      $z_{t-1}^{k-1}, M_t^{k-1} \leftarrow DM(z_t^{k-1}, P^{k-1}, t, \xi)$ ;  
7:      $M_t^k \leftarrow DM(z_t^k, P^k, t, \xi)$ ;  
8:      $\widehat{M}_t^k \leftarrow Edit(M_t^{k-1}, M_t^k, t)$ ;  
9:      $z_{t-1}^k \leftarrow DM(z_t^k, P^k, t, \xi) \{M \leftarrow \widehat{M}_t^k\}$ ;  
10:    if local then  
11:       $\alpha \leftarrow B(\overline{M}_{t, (w^{k-1})^k}) \cup B(\overline{M}_{t, (w^k)^k})$ ;  
12:       $z_{t-1}^k \leftarrow (1 - \alpha) \odot z_{t-1}^{k-1} + \alpha \odot z_{t-1}^k$ ;  
13:    end if  
14:  end for  
15: end for  
16:  $x^1, x^2, \dots, x^K = Decode(z_0^1, z_0^2, \dots, z_0^K)$ .
```

1.2. Training Details

We choose the pre-trained Stable Diffusion v1.4 in our pipeline. When sampling the original and the edited images, we run 50 inference steps with a classifier-free guidance scale of 7.5. During the optimization of strokes, we train 1000 iterations for each example. We use 96 strokes each of which includes 4 control points to represent the vector sketches. The

stroke width is defined with a fixed value 1.0. In the stroke-level local editing scheme, we adopt cross-attention maps of resolution 16×16 in up and down blocks in the diffusion model.

2. MORE RESULTS

We show more results including:

- Comparisons in **Word Swap** mode: Fig. 1;
- Comparisons in **Prompt Refinement** mode: Fig. 2;
- **Attention Re-weighting** mode: Fig. 3;
- **Iterative editing**: Fig. 4 and 5.

3. REFERENCES

- [1] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or, “Prompt-to-prompt image editing with cross-attention control,” in *ICLR*, 2023.

*The corresponding author is Chengying Gao.

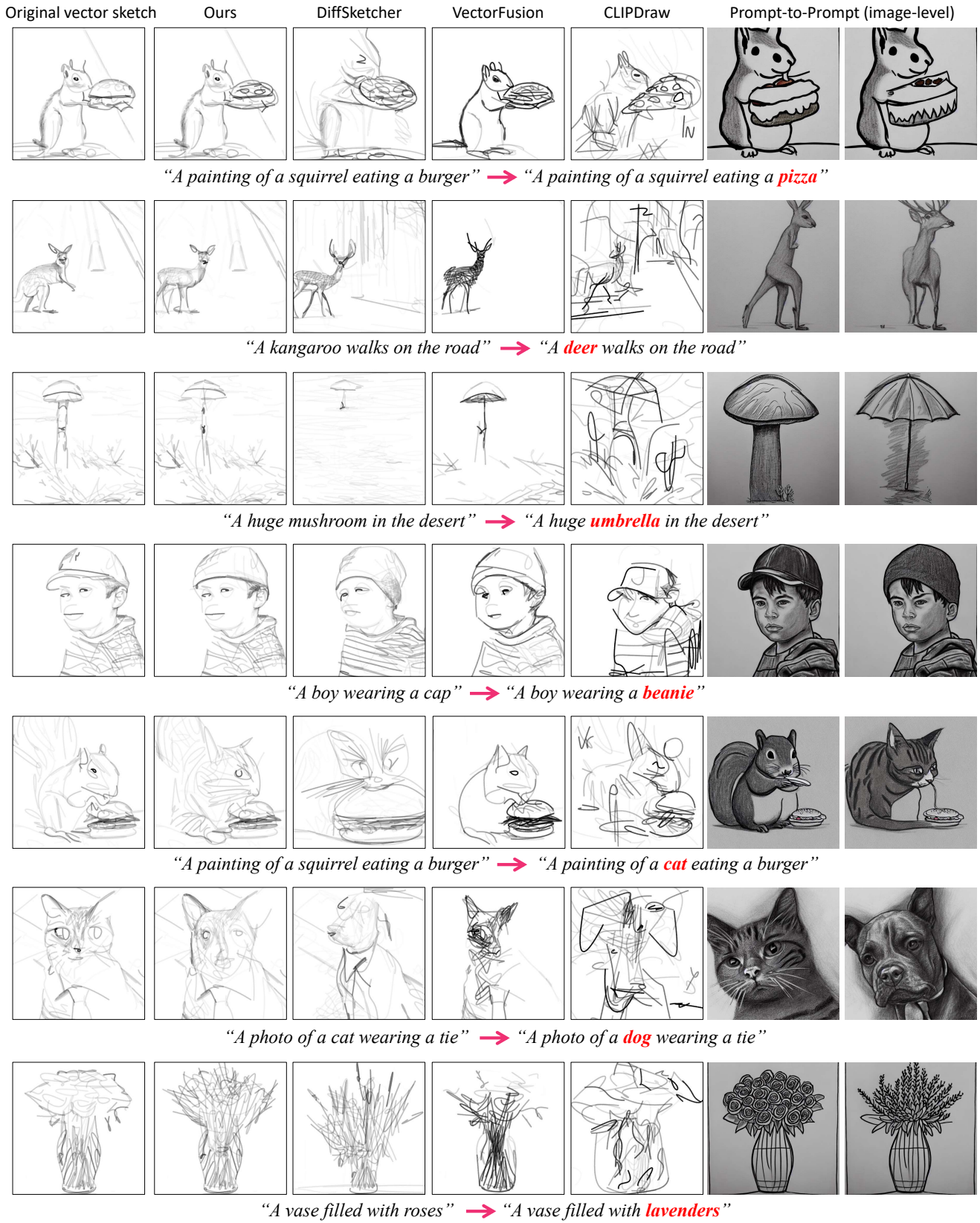


Fig. 1. Comparisons with baseline methods in Word Swap mode.

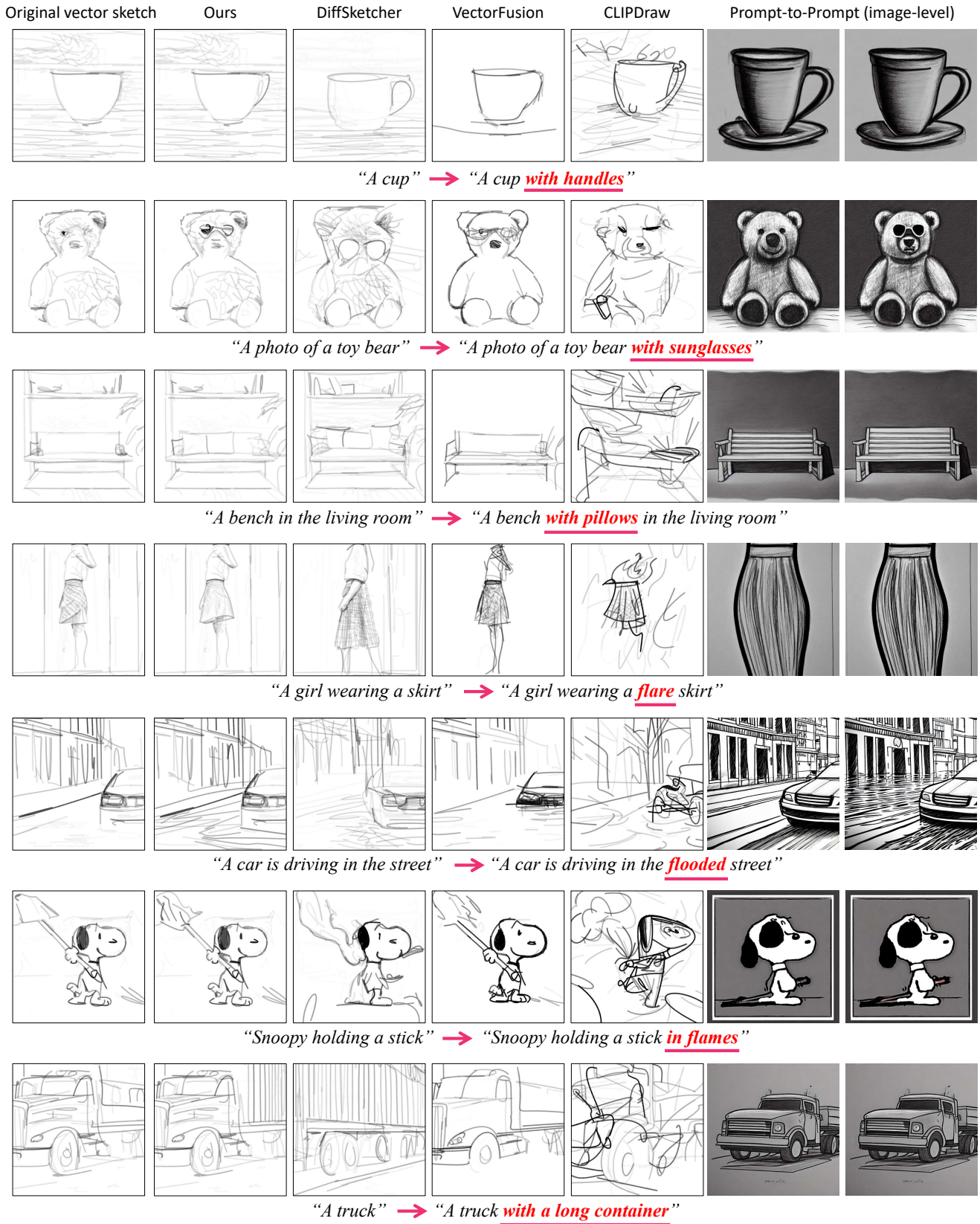


Fig. 2. Comparisons with baseline methods in Prompt Refinement mode.

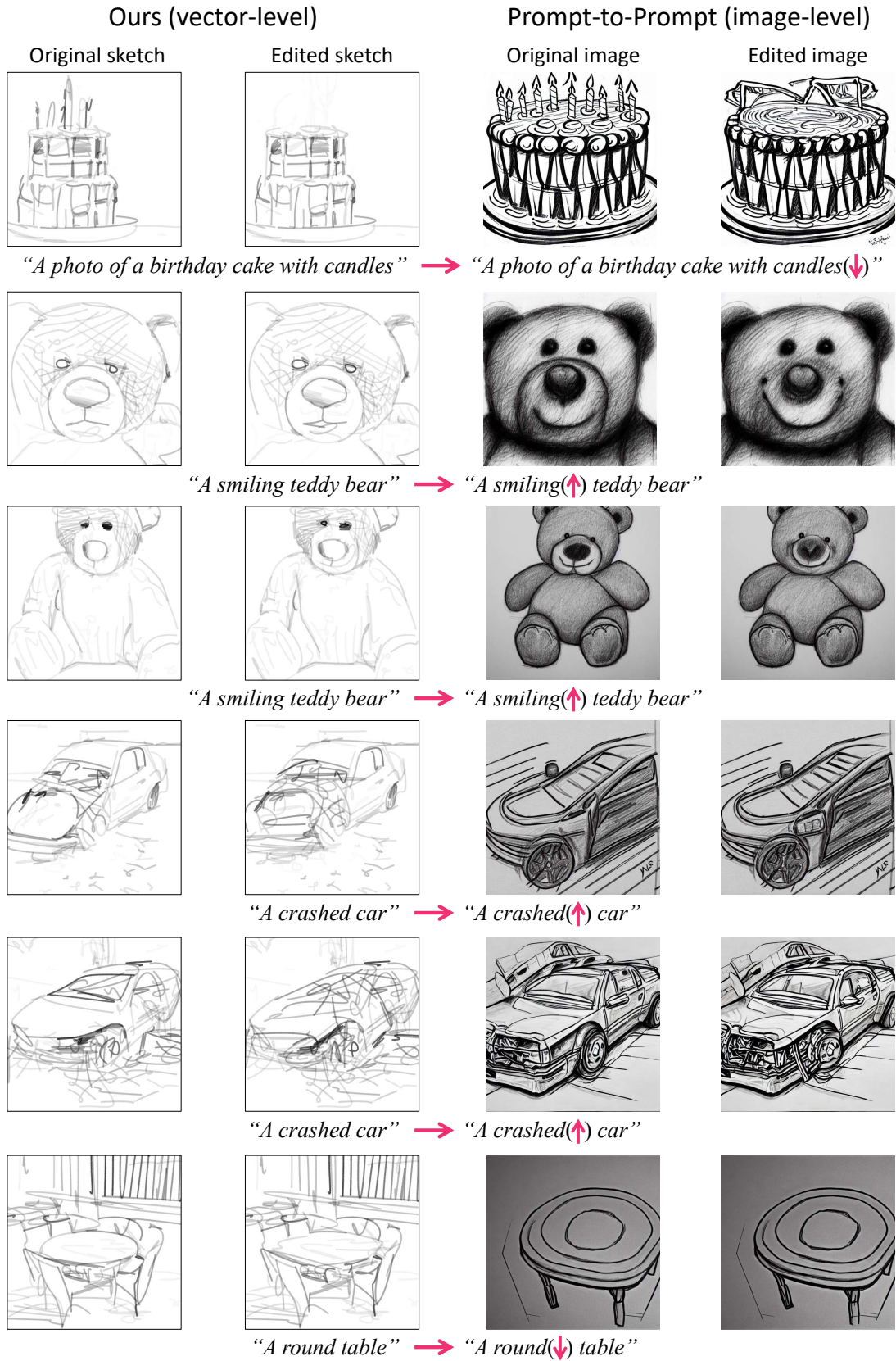


Fig. 3. Results in Attention Re-weighting mode.

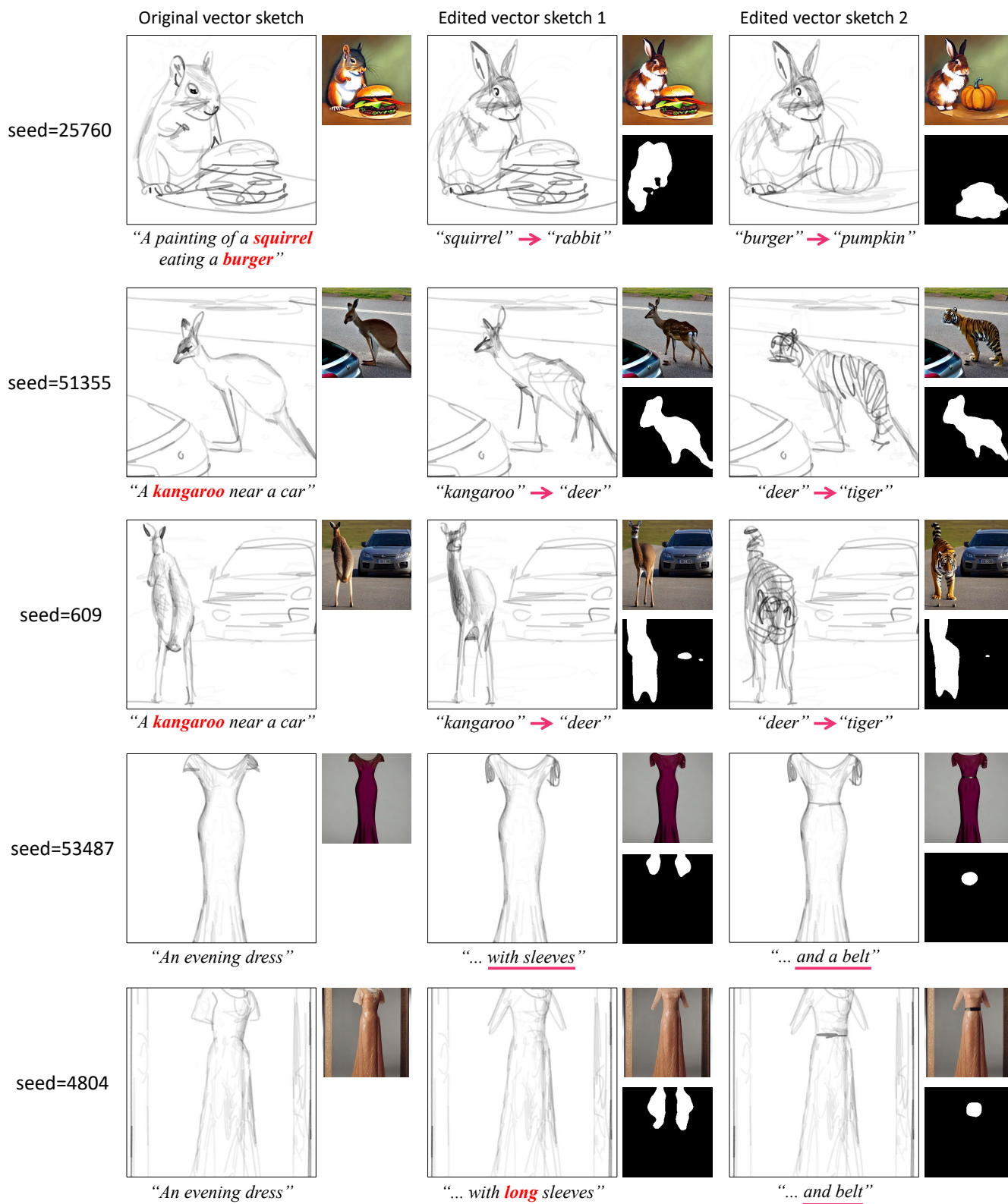


Fig. 4. Results of iterative editing.

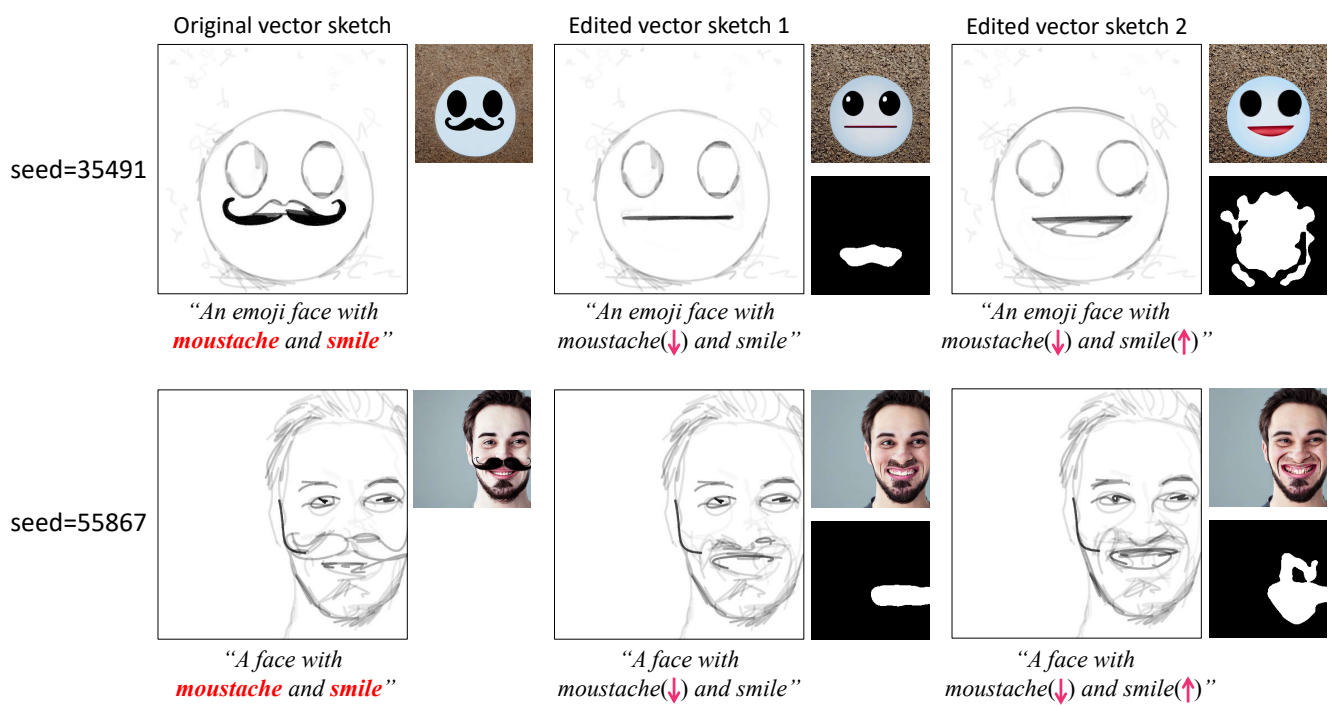


Fig. 5. Results of iterative editing.