

# Object-aware Inversion and Reassembly for Image Editing

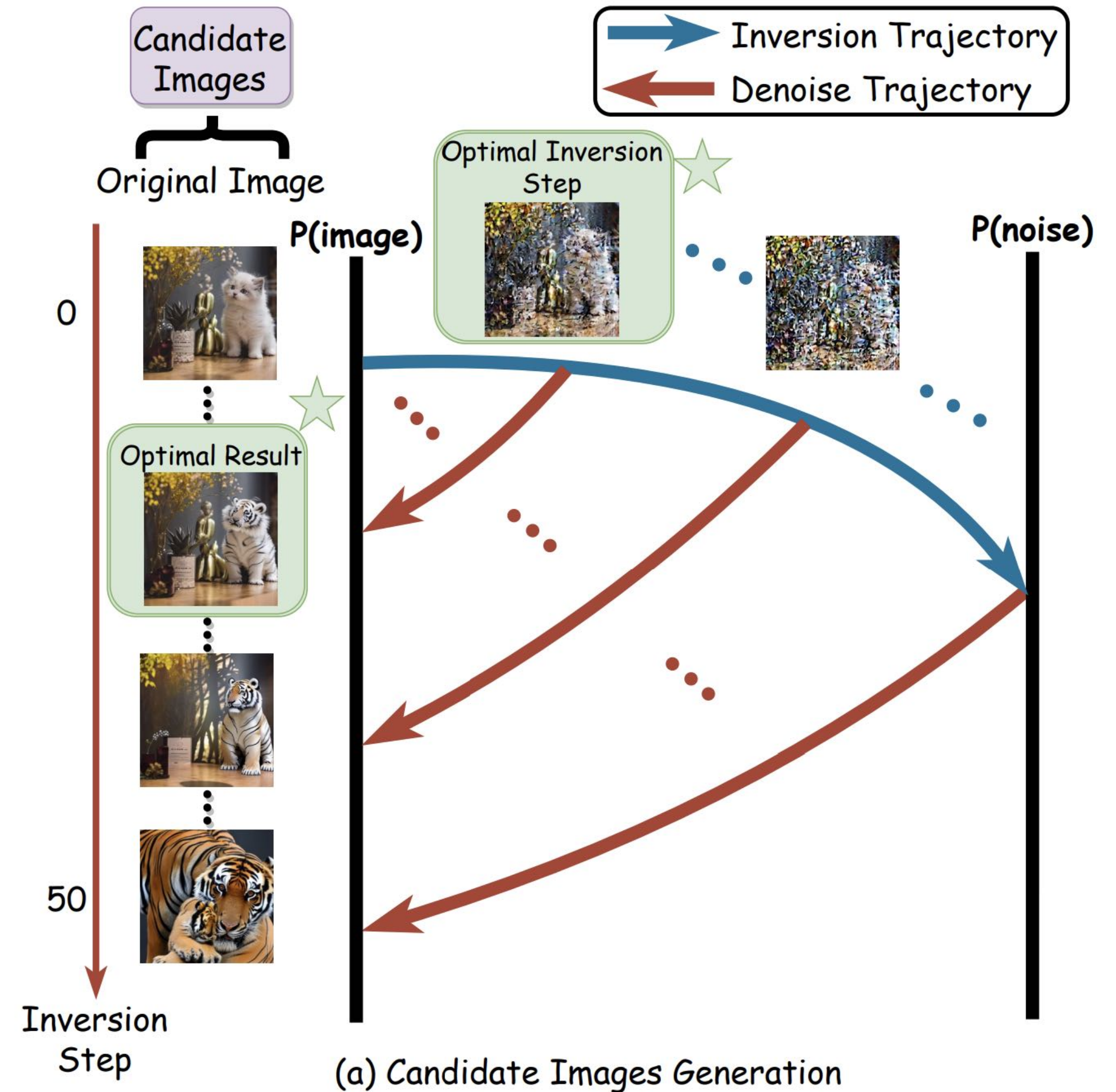
Zhen Yang<sup>1</sup>, Gangui Ding<sup>1</sup>, Wen Wang<sup>1</sup>, Hao Chen<sup>1</sup>, Bohan Zhuang<sup>2</sup>, Chunhua Shen<sup>1</sup>

1. Zhejiang University, 2. Monash University

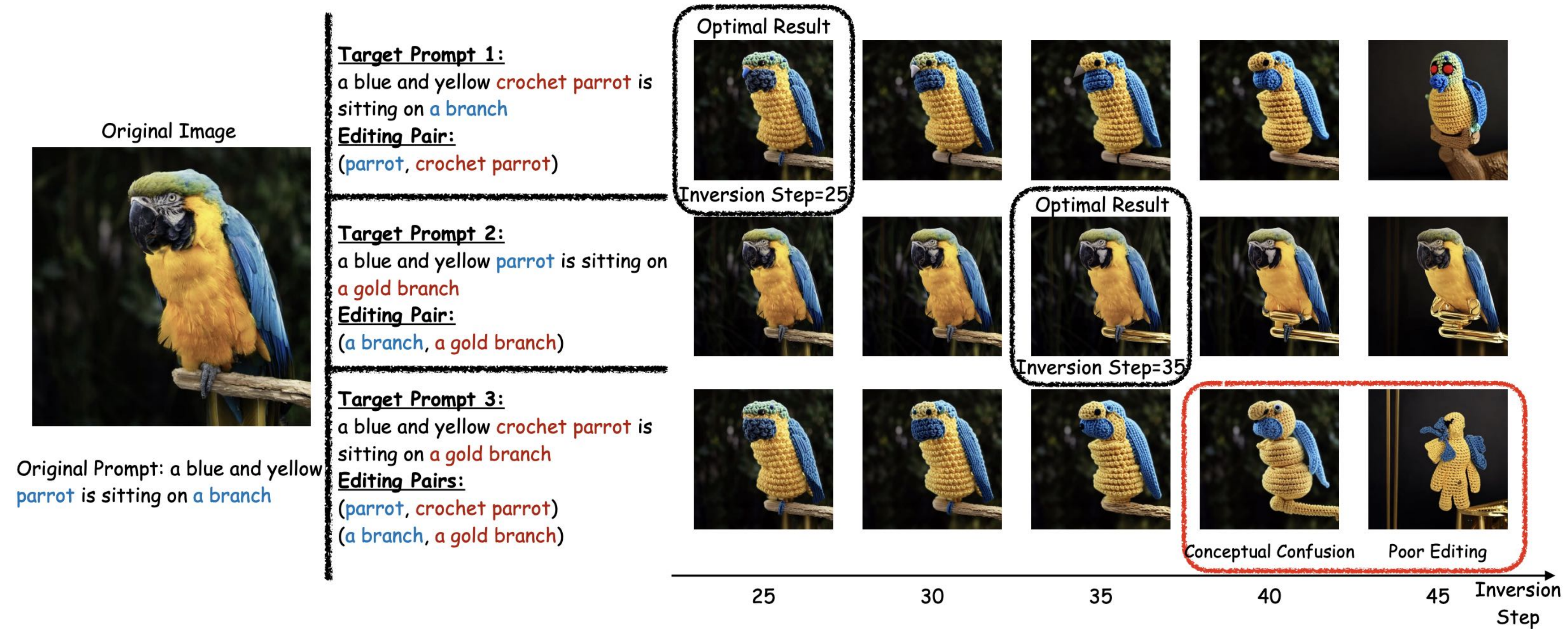


# Start Point

- Collect candidate images and manually select the optimal editing result.
- From  $t=0$  to  $t=50$ , the editing result ranges from "fidelity" to "editability".
- **People ignore "t".**
- Two technical problems: speed and manually. 🐒



# Motivation



- 1 Different editing pairs require different inversion steps !!! → Search Metric
- 2 Editing multiple pairs simultaneously can result in concept mismatch and poor editing !!! → OIR 🐒

# 1 Different editing pairs require different inversion steps !!

How to automatically select the optimal inversion step ? (Search metric)

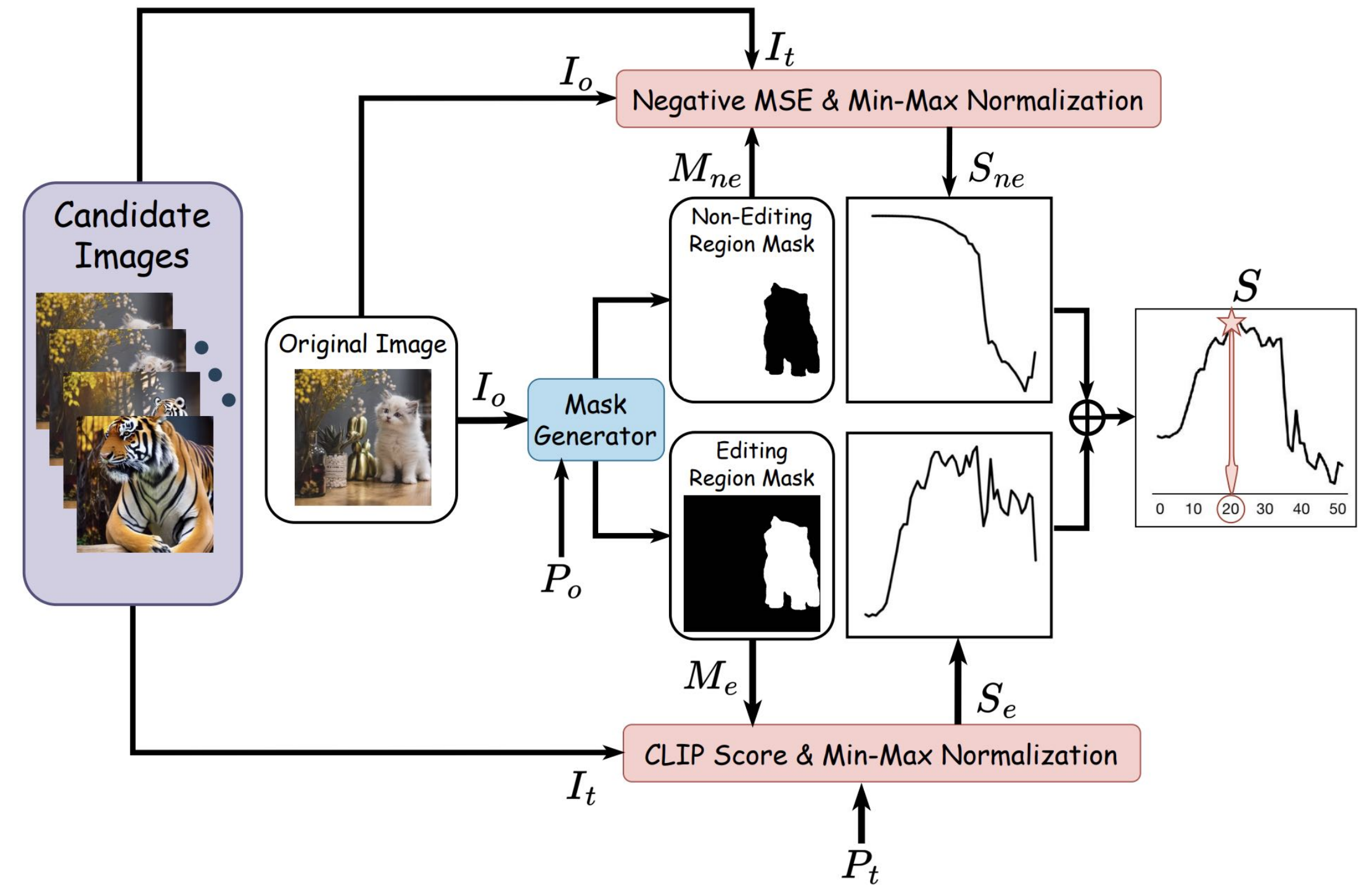
- Inspired by human observation, we hope that the **editing region will be aligned with the target prompt**, and the **non-editing region will be consistent to the original image**.
- Editing region metric: CLIP Score ( $\uparrow$ , 0→50)
- Non-editing region metric: Negative MSE ( $\downarrow$ , 0→50)

$$S_e(I_t, P_t, M_e) = \text{normalize} \left( \frac{CLIP_{image}(I_t, M_e) \cdot CLIP_{text}(P_t)}{\|CLIP_{image}(I_t, M_e)\|_2 \cdot \|CLIP_{text}(P_t)\|_2} \right)$$

$$S_{ne}(I_t, I_o, M_{ne}) = \text{normalize} \left( - \left\| (I_t - I_o) \odot M_{ne} \right\|_2^2 \right)$$

$$\text{Search Metric} = 0.5 \cdot (S_e + S_{ne})$$

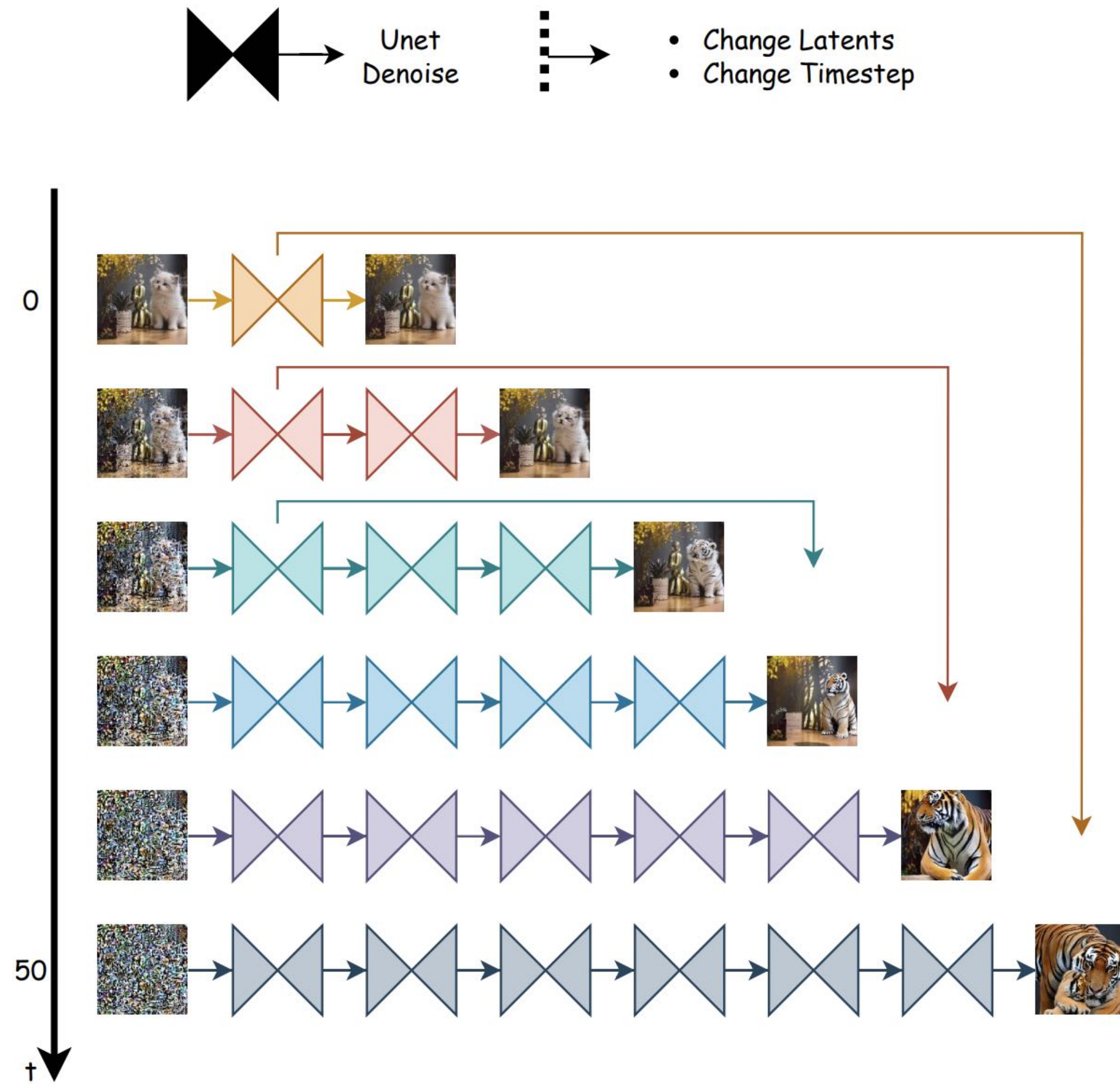
$P_o$ (original prompt): a **cat** is sitting on the desk } Editing Pair: (**cat**, **tiger**)  
 $P_t$ (target prompt): a **tiger** is sitting on the desk }



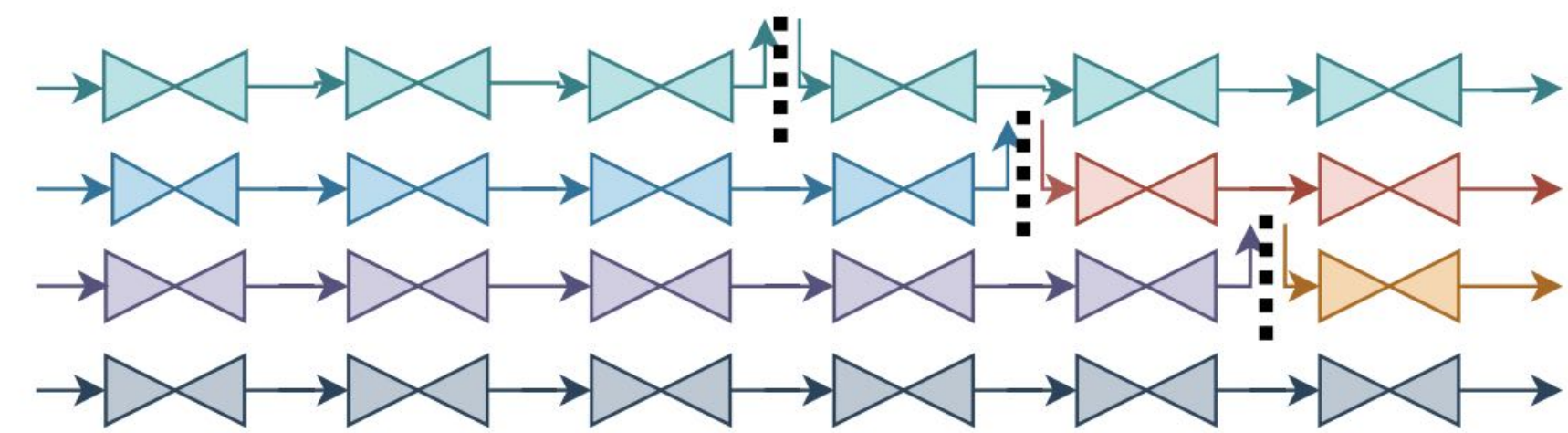
(b) Optimal Candidate Selection

# 1 Different editing pairs require different inversion steps !!

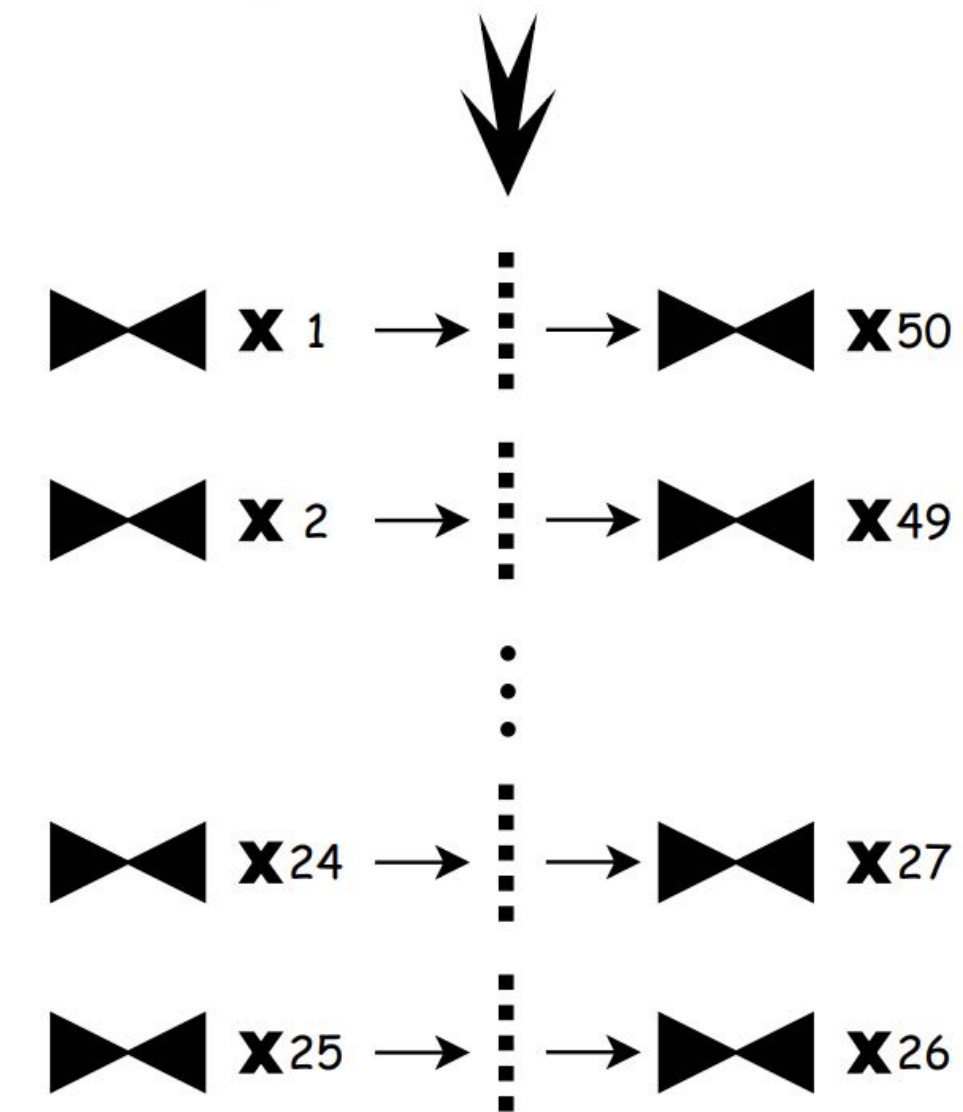
How to speed up ?



(a) Serial Denoise Methods



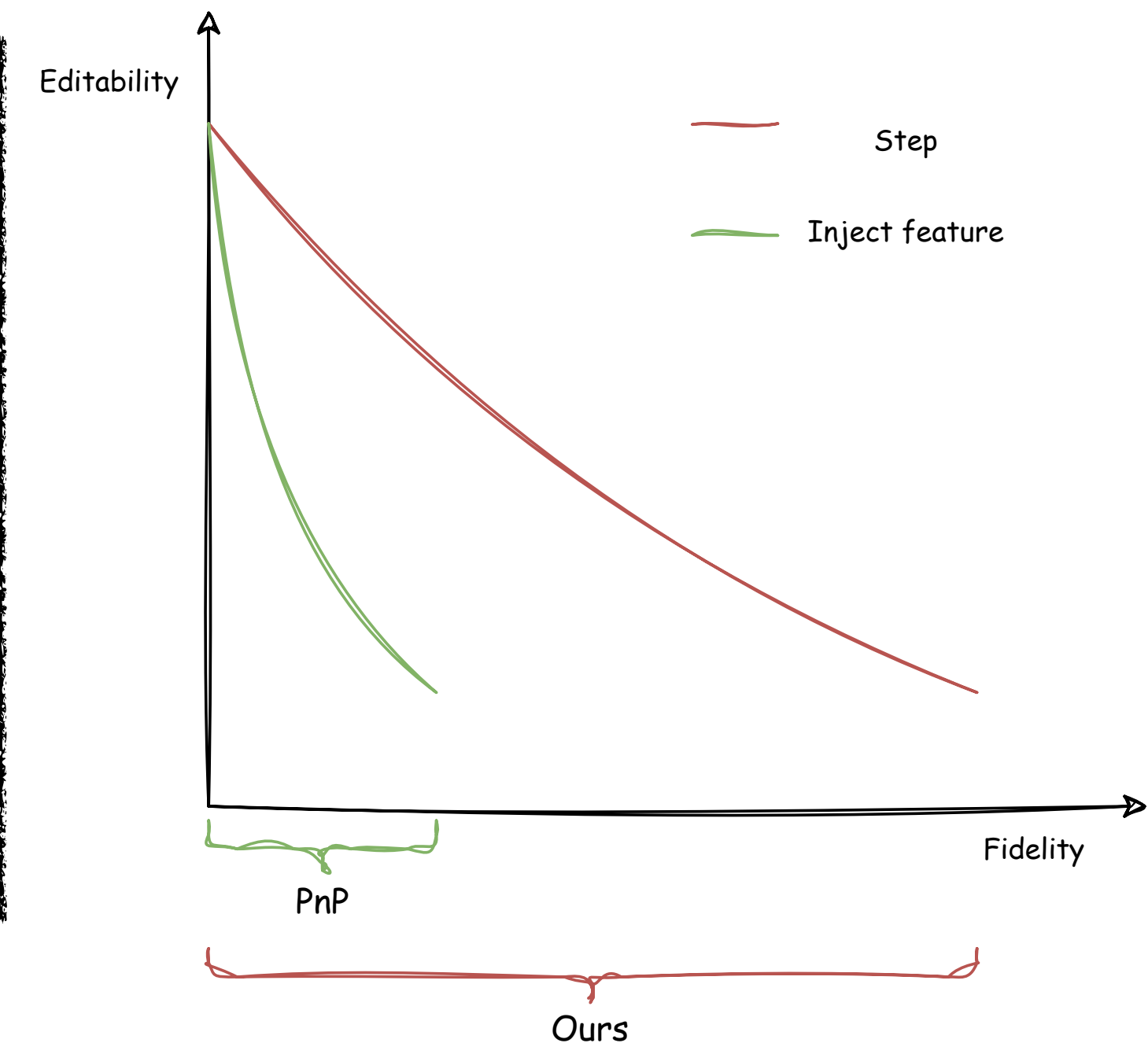
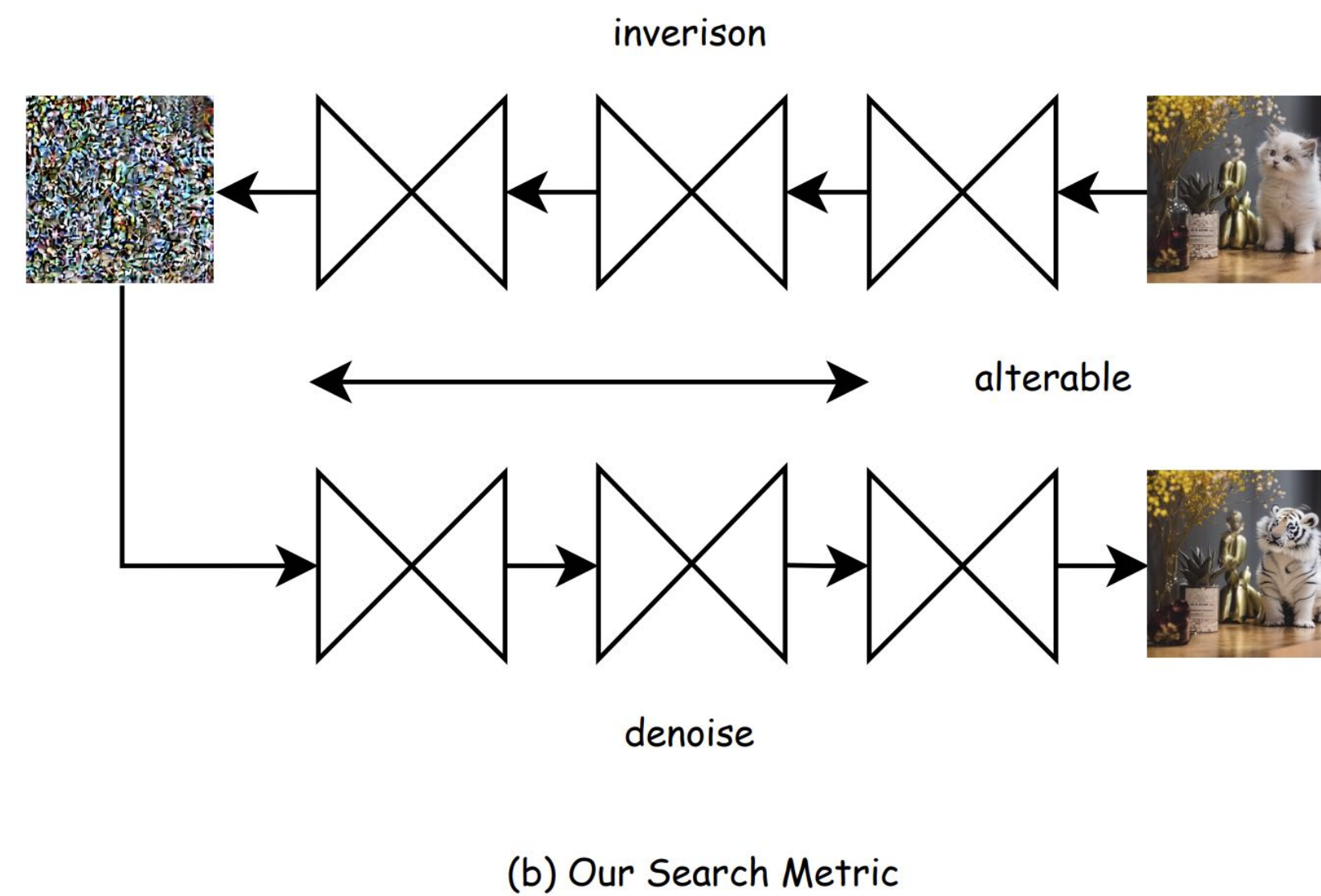
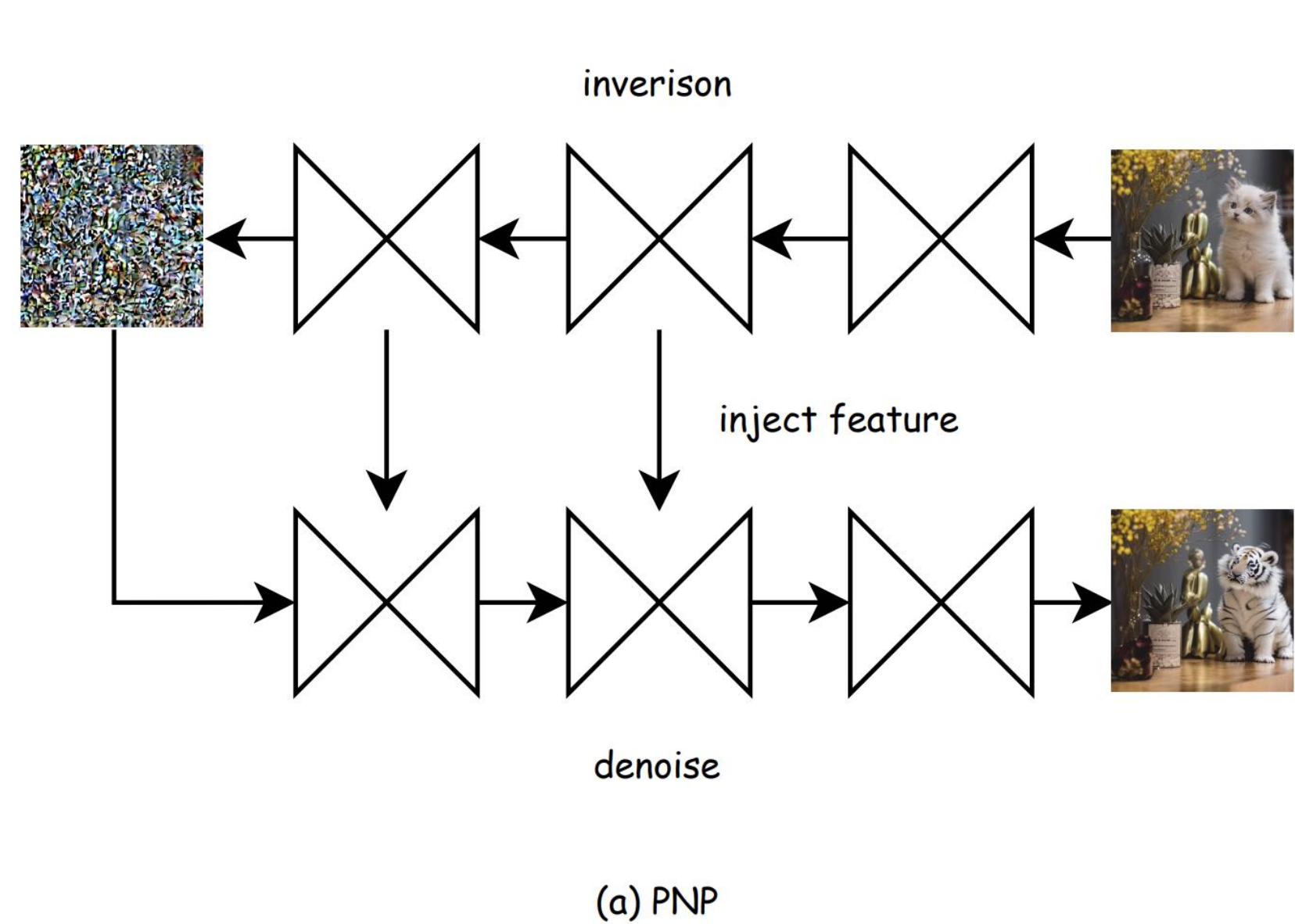
(b) Parallel Denoise Methods



(c) Parallel Denoise method Extends To 51-Step Denoise

# 1 Different editing pairs require different inversion steps !!

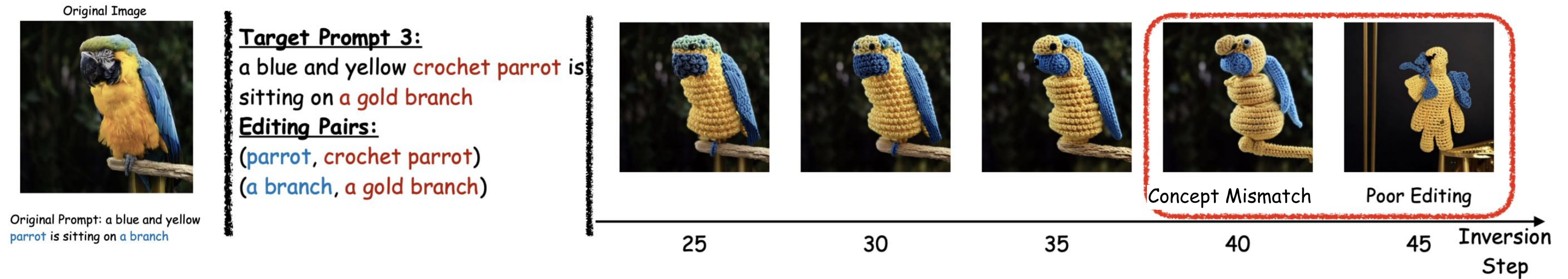
Why our novel method can work ?



- PNP injects features to control "editability" and "fidelity".
- Search metric changes the timestep to control "editability" and "fidelity".
- When PNP injects nothing, PNP is equivalent to the search metric with timestep=0. When PNP injects all features, PNP is equivalent to the search metric with timestep=50.

## 2 How to solve "concept mismatch" and "poor editing"!!!

What is the "concept mismatch" and "poor editing" and where are they from?



- **Concept mismatch:** mismatch the feature of **crochet** and **gold**.
  - Possible sources: There are overlapping features in the early stage of the denoising process.
- **Poor editing:** don't follow the objective of editing task (fidelity and editability).
  - Possible sources: Use the same inversion step for different editing pairs.

# 2 How to solve "concept mismatch" and "poor editing"!!!

How to solve these two problems ? (OIR)

Original Prompt( $P_o$ ): a dog sitting next to a gift

Target Prompt( $P_t$ ): ... cat ... cake ...

Editing Pairs: (dog, cat), (gift, cake)

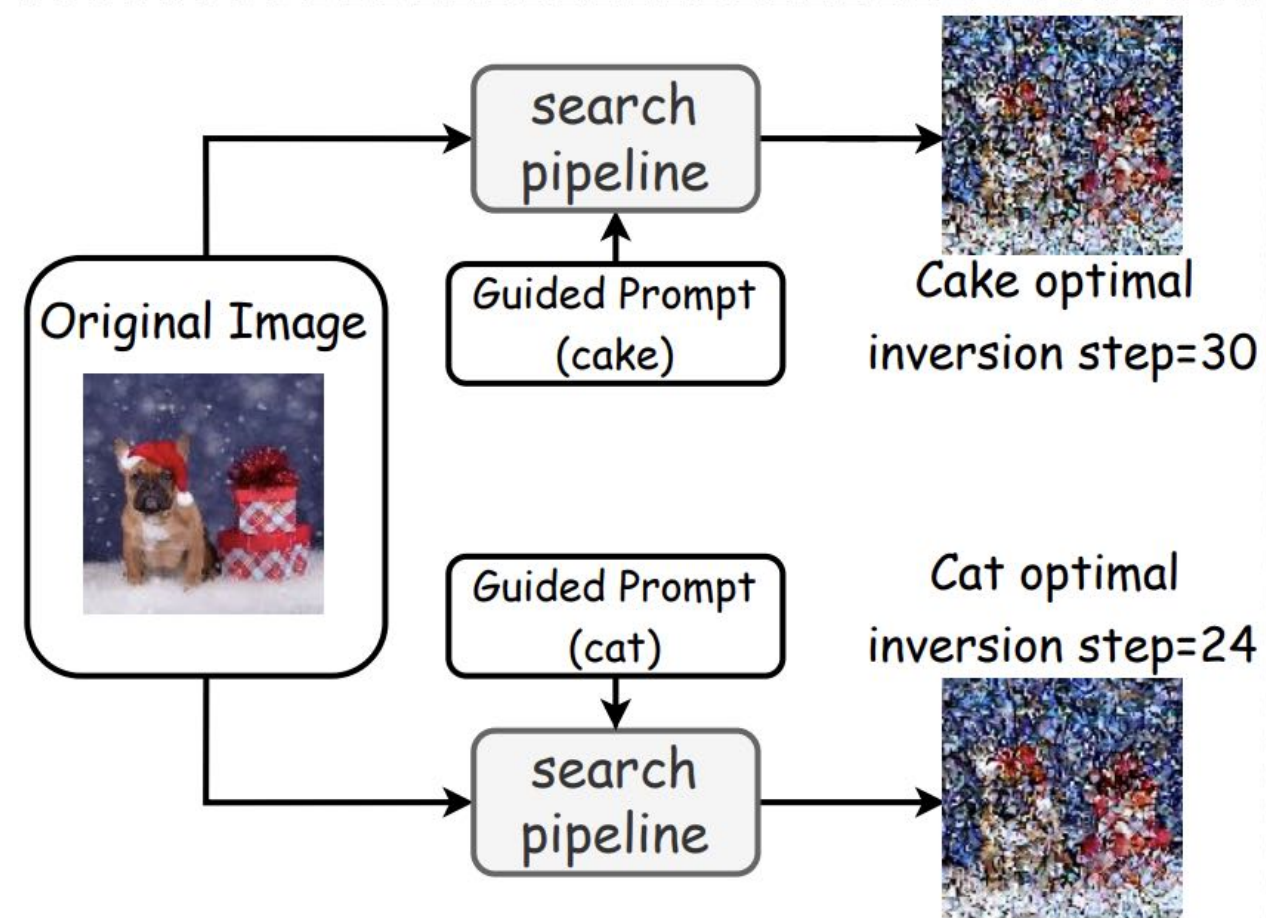
Guided Prompt(cat): ... cat ... gift ...

Editing Pair: (dog, cat)

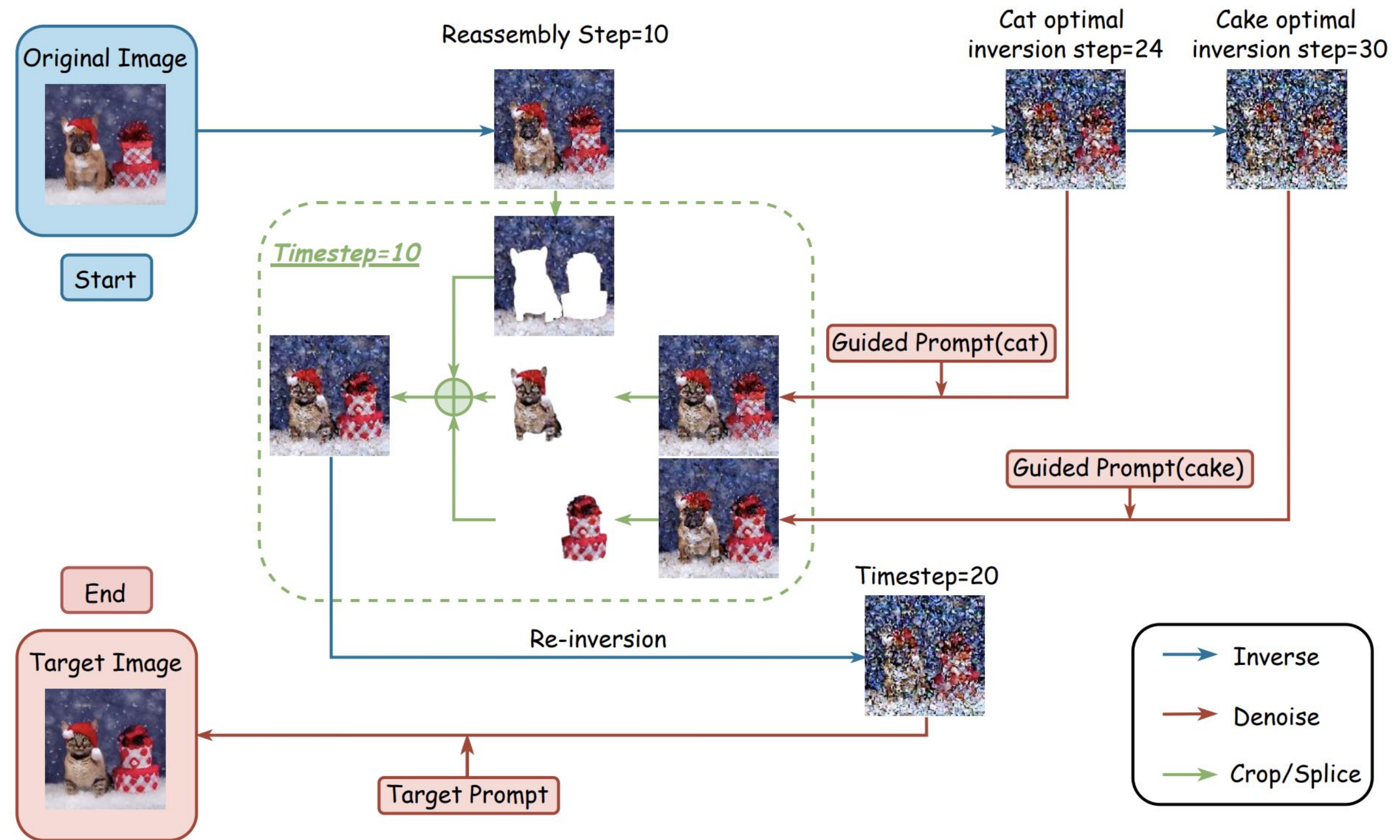
Guided Prompt(cake): ... dog ... cake ...

Editing Pair: (gift, cake)

(a) Guided prompts preparation



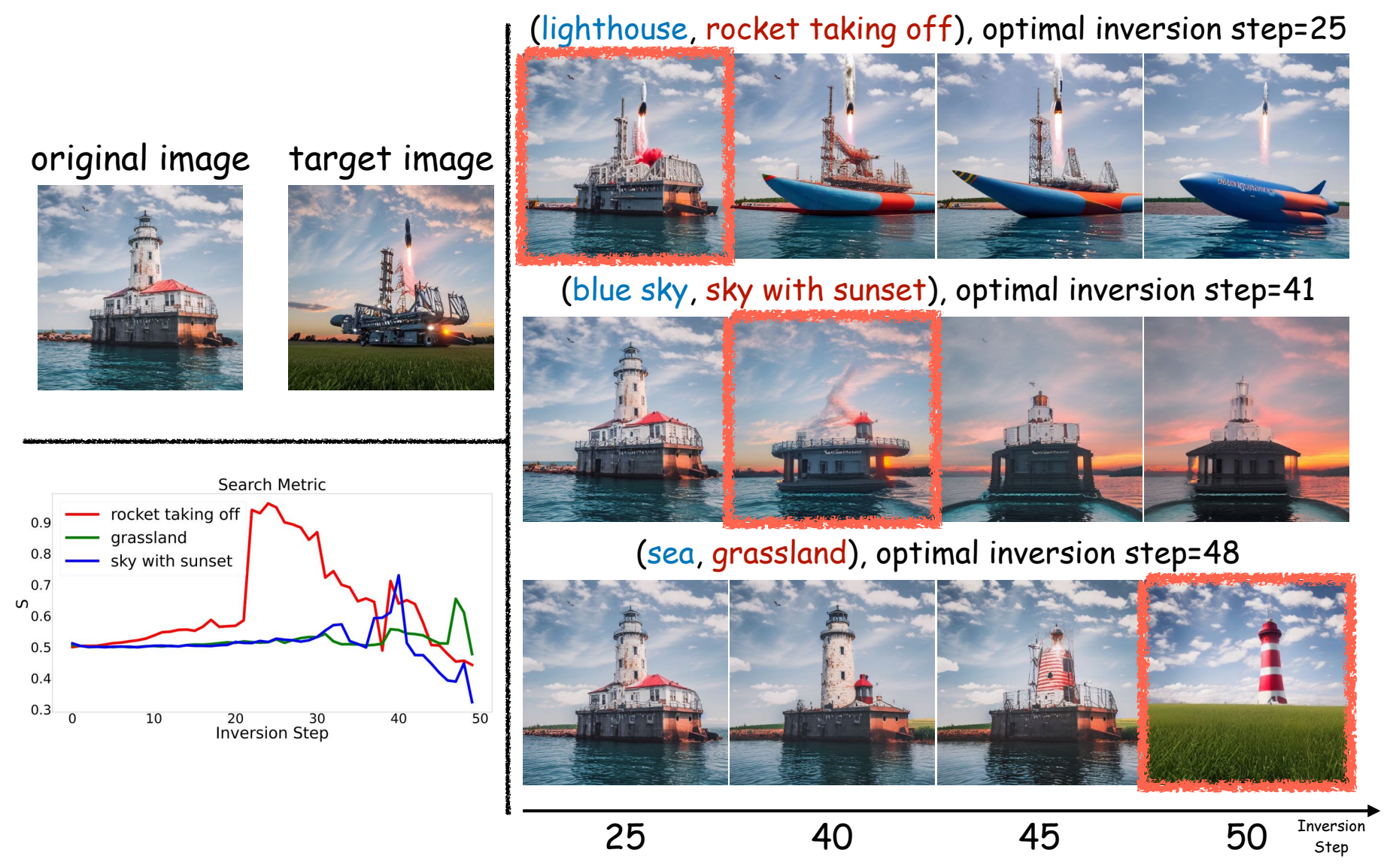
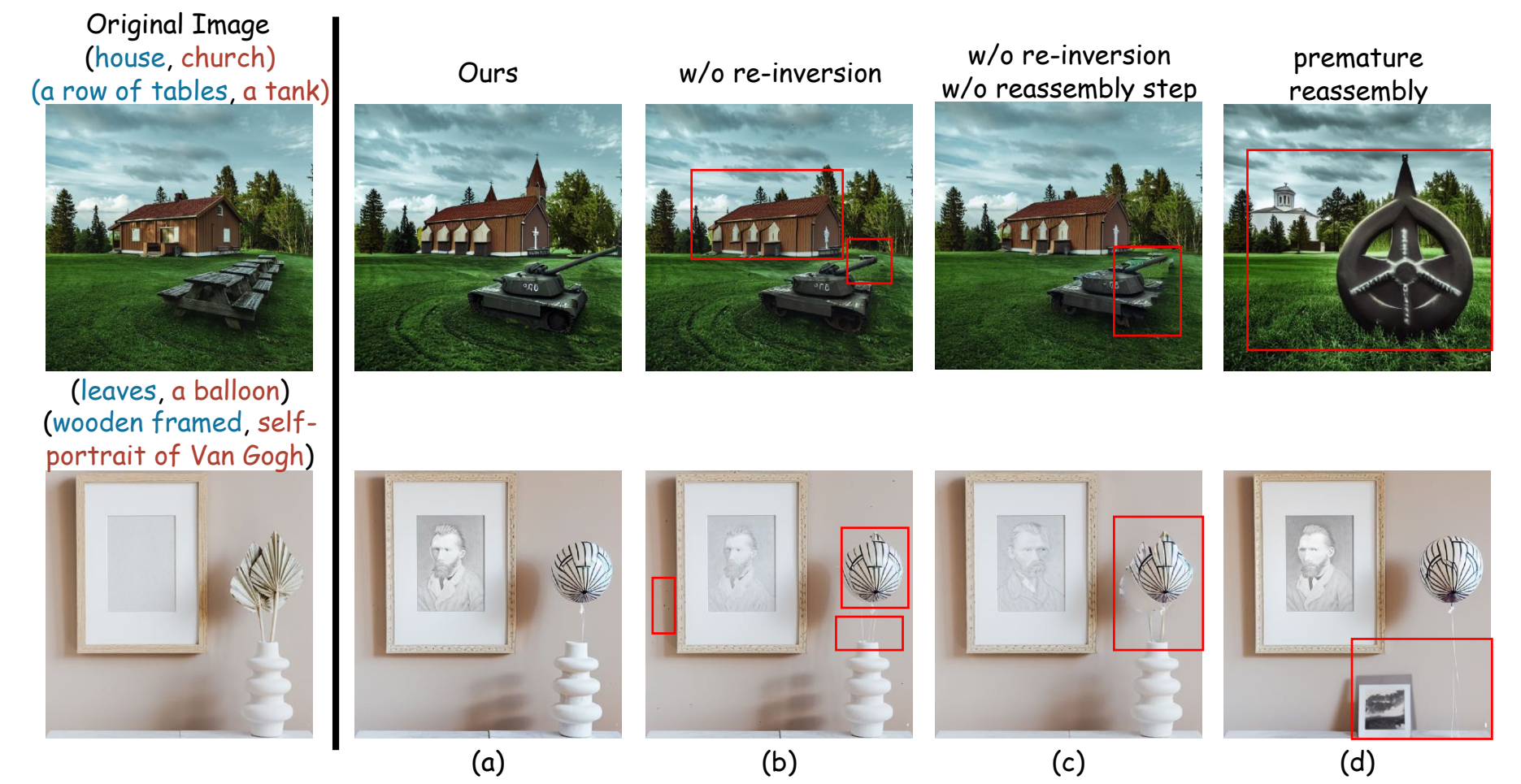
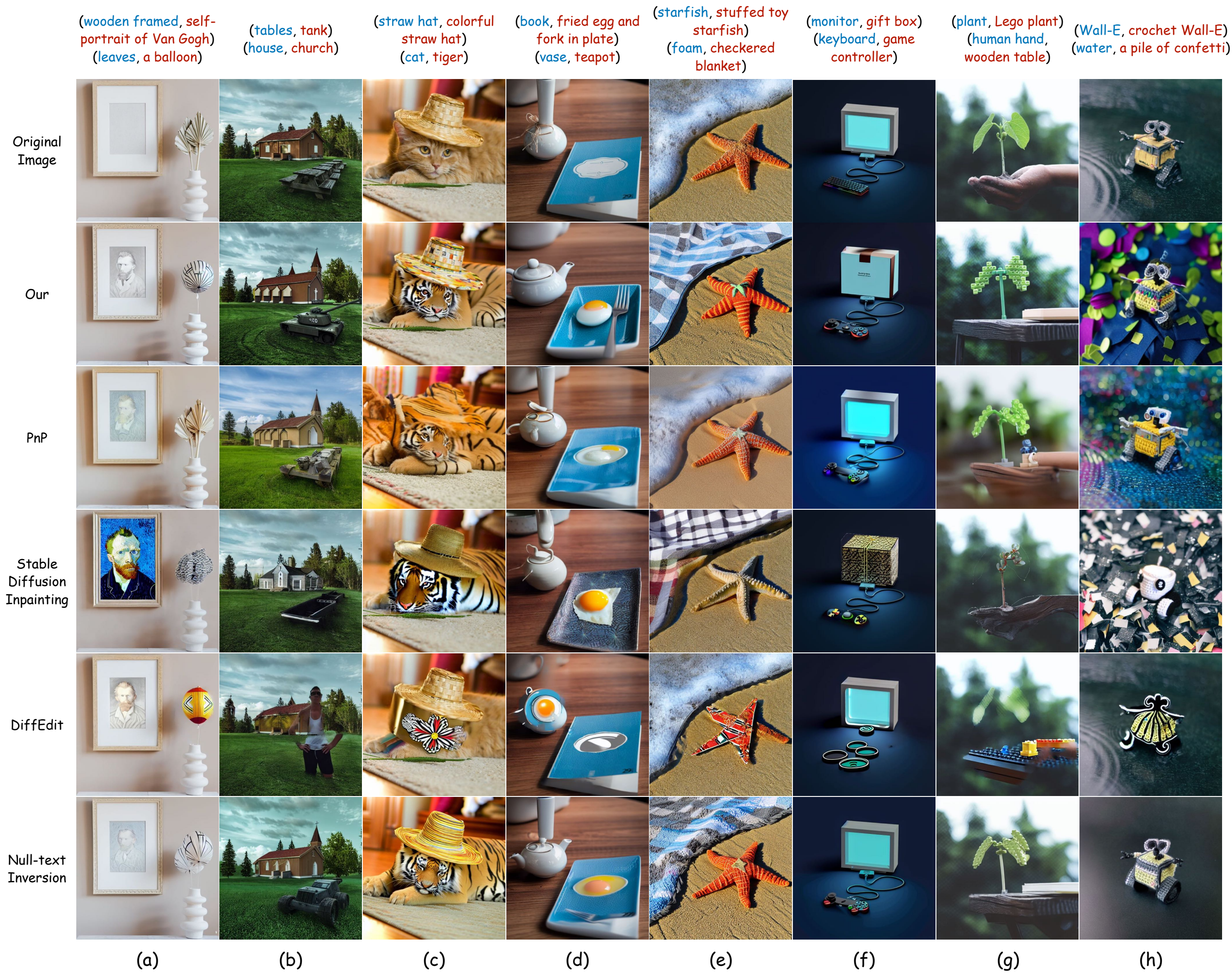
(b) Optimal inversion step searching



(c) Object-aware Inversion and Reassembly



# Results



I am currently seeking a Ph.D. supervisor.  
If you are interested, please contact me. 🤗



[zheny.cs@gmail.com](mailto:zheny.cs@gmail.com)