

# IntentEdit: Multi-Agent Reasoning for Intent-Driven Complex Image Editing

## Supplementary Material

### A. More Implementation Details

**RISE Bench.** RISEBench [2] is a recently proposed benchmark designed to evaluate reasoning-informed visual editing capabilities of generative models. Unlike prior benchmarks that mainly test superficial modifications (e.g., color changes or object replacement), RISEBench systematically focuses on four fundamental dimensions of human-like reasoning in visual editing: *Temporal*, *Causal*, *Spatial*, and *Logical* reasoning. Each dimension is carefully constructed with diverse, high-quality test cases that require models to go beyond pattern matching and demonstrate reasoning-driven visual understanding. For example, temporal reasoning evaluates whether models can capture natural progressions (e.g., fruit ripening), causal reasoning requires models to reflect effects of external forces (e.g., collisions), spatial reasoning involves structural coherence and viewpoint transformations, and logical reasoning tests puzzle-solving and rule-based inference (e.g., tic-tac-toe).

In our experiments, we strictly follow the official implementation and default settings of RISEBench, generating images and conducting evaluations exactly as specified by its authors, thereby ensuring fair and directly comparable results.

**Complex-Edit Bench.** Complex-Edit [1] is a complexity-controllable benchmark designed to evaluate instruction-based image editing under varying levels of compositional difficulty. Unlike benchmarks with atomic or low-level edits, Complex-Edit constructs instructions in a *chain-of-edit* manner: starting from a set of atomic editing operations (e.g., adding an object, adjusting color, changing background), GPT-4o is employed to generate sequential instructions, which are then simplified and compounded into increasingly complex editing directives. This allows fine-grained control over instruction complexity, ranging from simple single-step edits (C1) to highly intricate transformations with multiple dependencies (C8).

The benchmark evaluates models across three primary dimensions: *Instruction Following*, *Identity Preservation*, and *Perceptual Quality*. A VLM-based auto-evaluation pipeline is adopted to ensure scalable and reproducible assessment. In our implementation, we specifically focus on the subset of real images provided by Complex-Edit. We evaluate all models at complexity setting (C5), and each image is assessed with five independent evaluation runs to mitigate per-sample variance. This setup ensures both robustness and comparability with prior work, while also highlighting model behavior under realistic, challenging editing scenarios.

### B. Samples of RISE-MIX Bench

Figure 1 presents representative examples from our RISE-MIX benchmark. Unlike existing datasets that primarily focus on either atomic edits or synthetically composed instructions, RISE-MIX emphasizes realistic scenarios where reasoning and direct edits are inherently intertwined. Each sample requires the model to not only interpret abstract reasoning cues (e.g., imagining object growth, predicting plausible temporal changes) but also to execute precise, visually grounded modifications (e.g., object replacement, background alteration). This dual nature makes RISE-MIX particularly suited for evaluating the capability of multimodal systems to bridge high-level reasoning with low-level visual editing.

### C. Generation Process Visualization

To provide a more intuitive understanding of our framework, we present several qualitative examples that illustrate the full generation process. Each case demonstrates how the model interprets the instruction, executes editing steps, and then reflects the output.

As shown in Figure 2 and Figure 3, our framework first analyzes and reasons about the user intent, and then decomposes a complex instruction into multiple intermediate steps. At each step, the model performs targeted editing operations (e.g., duplicating objects, adjusting colors, or modifying spatial arrangements) and subsequently generates a reflection assessing the success of the edit. Reflections are categorized into Success, Refinement, or Replan, which guide the model to either proceed, refine the current output, or retry the step. This step-wise Plan-Edit-Reflect process enables the system to self-monitor its progress and ensures higher fidelity to the original instruction.

### D. Additional Qualitative Comparison

As shown in Figures 4 and 5, we present additional qualitative results from the RISE Bench and Complex-Edit Bench. Additionally, we present a qualitative comparison of our two variants against several advanced closed-source models, including GPT-Image, Nano Banana, and SeedEdit3.

### E. User Study

To complement the automatic evaluations provided by RISEBench, RISE-MIX Bench and Complex-Edit Bench, we conducted a user study to further assess the perceived quality of model outputs from a human perspective. Our

goal was to examine whether the automatic metrics align with human judgments and to provide a more comprehensive understanding of model performance in realistic usage scenarios.

Table 1. **User study.** Results of human evaluation for complex image editing. Participants assessed model outputs on three criteria: Identity Preservation (IP), Instruction Following (IF), and Visual Quality (VQ). Boldface denotes the best performance, while underlining indicates the second best.

Method	IP $\uparrow$	IF $\uparrow$	VQ $\uparrow$	Overall $\uparrow$
Step1X-Edit	3.4	3.3	3.5	3.4
ICEdit	3.6	2.9	3.5	3.3
Ovis-U1	3.7	3.9	4.1	3.9
OmniGen2	3.6	3.8	3.9	3.8
IEAP	2.1	2.5	2.3	2.3
CCA	1.9	2.1	2.6	2.2
BAGEL	3.3	4.0	3.8	3.7
BAGEL-CoT	3.5	3.9	3.1	3.5
Kontext-dev	4.0	3.5	<b>4.5</b>	4.0
Ours(InternVL3)	<b>4.2</b>	<u>4.1</u>	<u>4.1</u>	<u>4.1</u>
Ours(GPT-4o)	<u>4.1</u>	<b>4.4</b>	<u>4.2</u>	<b>4.2</b>

## E.1. Setup

We recruited 10 participants with backgrounds spanning computer science, design, and non-technical fields. All participants reported normal or corrected-to-normal vision. Each participant was presented with pairs of input-output images generated by different models, accompanied by the corresponding editing instruction. To avoid bias, images were shuffled and anonymized, and participants were not informed of which model produced which output. Each participant evaluated 26 randomly sampled editing tasks drawn from both RISE bench, Complex-Edit bench and RISE-MIX bench. Each task were generated from a combination of our method and 9 baseline algorithms. Specifically, the tasks consisted of randomly sampled editing tasks drawn from three distinct benchmark datasets: RISE Bench, Complex-Edit Bench, and RISE-MIX Bench. Participants were asked to rate the model outputs along three dimensions using a 5-point scale:

1. **Instruction Following(IF):** How well does the output reflect the specified instruction?
2. **Identity Preservation(IP):** To what extent are irrelevant aspects of the input image preserved?
3. **Visual Quality(VQ):** How realistic, coherent, and aesthetically pleasing is the output?

In total, the study yielded a comprehensive dataset consisting of 26 tasks evaluated by 10 participants, resulting in 11 ratings per task. This amounted to a total of  $26 \times 10 \times 11$

$\times 3 = 8580$  ratings, providing a robust set of evaluations for further analysis.

## E.2. Results

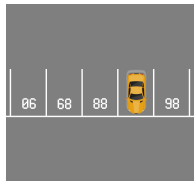
The user study results, as shown in Table 1, demonstrate the superior performance of our proposed system. The GPT-4o variant achieves the highest overall score of 4.2. In comparison, the InternVL3 version scores 4.1 overall, ranking second. Ovis-U1 and Kontext-dev show competitive results, with Kontext-dev achieving the best VQ score of 4.5, while Ovis-U1 ranks high in IF and VQ. BAGEL and BAGEL-CoT perform well in IF but do not surpass our models in overall performance. Baseline methods such as Step1X-Edit, ICEdit, and IEAP score significantly lower, particularly in IF and IP, highlighting their limitations in handling complex image editing tasks. These results underline the effectiveness of our approach in intent-driven complex image editing tasks.

## F. Limitation

Our multi-agent editing framework inevitably introduces a trade-off between multi-step capability and image quality: successive multi-turn edits may accumulate artifacts, even though we partially mitigate this with restoration. In addition, iterative planning increases inference time compared to single-pass approaches. Future work will focus on exploring acceleration techniques for diffusion models and designing diffusion architectures that better preserve image quality under multi-step editing.

## References

- [1] Siwei Yang, Mude Hui, Bingchen Zhao, Yuyin Zhou, Nataniel Ruiz, and Cihang Xie. Complex-edit: Cot-like instruction generation for complexity-controllable image editing benchmark. [arXiv preprint arXiv:2504.13143](#), 2025. 1
- [2] Xiangyu Zhao, Peiyuan Zhang, Kexian Tang, Xiaorong Zhu, Hao Li, Wenhao Chai, Zicheng Zhang, Renqiu Xia, Guangtao Zhai, Junchi Yan, et al. Envisioning beyond the pixels: Benchmarking reasoning-informed visual editing. [arXiv preprint arXiv:2504.02826](#), 2025. 1



Draw what it will look like after the car moves out from the parking slot and add a 'Reserved' sign to the parking slot where the car was previously parked. Additionally, change the color of the background to a lighter shade of gray.



Draw what it will look like after being sewn together and add decorative stitching along all edges of the sewn fabric. Then place the sewn fabric onto a rustic wooden table as the background.



Draw what it will look like during the next king tide, add a beach umbrella in the foreground, and place a lifeguard tower near the left side of the beach.



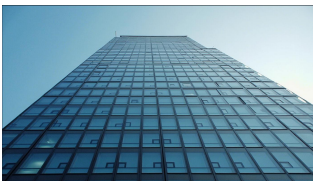
Draw what it looks like after being violently shaken and then suddenly opened. Additionally, replace the setting with a festive party environment and add colorful confetti flying in the air.



Draw what it will look like after a week in a damp basement and add scattered crumbs around the plate along with a faint watermark on the bread that says 'For Display Only.'



Draw what it will look like 50 years after abandonment, add graffiti art to some parts of the carousel, and replace the background with an overgrown forest.



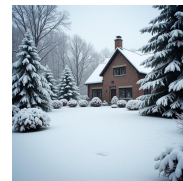
Draw the front view of this building as seen from a distance, add a fountain in front of the building, and create a cloudy sky in the background.



Draw what the cutting board will look like after three months of daily use and add knife marks on its surface. Also, place a small bowl of chopped vegetables on top of it.



Draw the left view of the objects in the image, and add a label that says 'Cube' on the cube and 'Sphere' on the sphere. Additionally, change the background color to light blue.



Draw what it will look like after many people have walked on it, add a snowman near the front door, and place a wooden bench with snow piled on it under one of the pine trees.



Draw what it will look like in spring, with leaves fully grown on the tree branches, and add blooming flowers around the base of the tree along with a wooden bench placed nearby.



Draw the computer after the side cover has been removed, add a glowing LED light strip inside the case that changes colors, and place the computer on a wooden desk with a mouse pad and a wireless mouse next to it.



Draw what it will look like after being deep-fried, add a dusting of chopped parsley on top, and place it on a rustic wooden platter.



Draw what it will look like after a century of global warming, add a large abandoned signpost with 'Restricted Area' written on it near the glacier, and change the sky to an ominous orange hue.



Draw what the keyboard will look like after five years of continuous use, add a coffee stain on one corner of the keyboard, and replace the background with a wooden desk.



Draw what it will look like after it is stepped on by muddy boots and add a small potted plant in the corner near the chair. Also, place a pair of neatly arranged slippers next to the ottoman.



Draw what it will look like after being splashed with ink. Additionally, change the color of the wooden table to a dark mahogany and fold the napkin into an origami rose shape.



Generate an image assembling the provided disassembled clock components into a complete clock displaying the time as 9:45. Additionally, replace the clock's metallic frame with a wooden frame and add a small engraving on the wooden frame that reads 'Timeless Elegance'.



Draw what the building will look like in twenty years, add solar panels to the roof, and replace the tree on the left with a fountain.



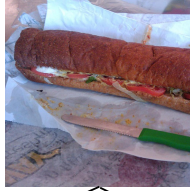
Draw what it will look like immediately after one of its edges is lightly burned, add a small handwritten note in the top-left corner saying 'Fragile', and place a subtle coffee stain in the bottom-right corner.



Draw what it looked like during the Jurassic period, and add a lush prehistoric jungle as the background with tall cycads and ferns. Include a flying Pterosaur gliding through the sky above.

Figure 1. Samples from the proposed **RISE-MIX Bench**, illustrating real-world, reasoning-intensive editing instructions that combine both abstract reasoning and concrete visual modifications. Each example pairs an input image with complex instructions and the corresponding target edits, showcasing diverse reasoning types.

Instruction: Place two identical sandwiches side-by-side on a rustic wooden table setting. Add a steaming cup of coffee in the background and adjust the lighting to warm tones with a soft ambiance. Change the color of the sandwich wraps to a deeper brown to match the aesthetic.



PLAN1

PLAN2

```
'steps': [{'step_number': 1,
'step_instruction': 'Place two identical sandwiches side-by-side on a rustic wooden table setting. Add a steaming cup of coffee in the background and adjust the lighting to warm tones with a soft ambiance. Change the color of the sandwich wraps to a deeper brown to match the aesthetic.',
'instruction_type': 'Direct',
'model_choice': 'Kontext'}]
```

STEP1



Reflect

"reason": "The core instruction of placing two identical sandwiches side-by-side with a rustic wooden table setting and deeper brown wraps was partially fulfilled, but the final image radically diverged. The sandwiches appear as round buns rather than matching those in the Source Image. Additionally, the wraps were not visibly adjusted to a deeper brown, as intended, and the 'rustic wooden table setting' lacks notable texture or aesthetic rusticity."

"quality score": 5

"reflection": "Replan" ❌

```
'steps': [{'step_number': 1, 'step_instruction': 'Duplicate the long baguette-style sandwich from the original image to create two identical sandwiches side-by-side with identical fillings, bread shape, and overall appearance.',
'instruction_type': 'Reasoning', 'model_choice': 'Bagel'}, {'step_number': 2, 'step_instruction': 'Place the sandwiches on a rustic wooden table with visible aged grain patterns and warm brown tones. Add a steaming cup of coffee in the background, with steam clearly visible to emphasize freshness and ambiance. Change the sandwich wrapping to a darker brown color for a cohesive rustic aesthetic. Adjust the lighting to warm, soft tones that evoke a cozy ambiance.',
'instruction_type': 'direct',
'model_choice': 'Kontext'}]
```

STEP1



STEP2



Reflect

"quality score": 5

"reflection": "Success" ✅

Instruction: Draw the appearance of the image after the camera has moved to a position very close to the train, add a graffiti art with colorful designs on the side of the train, and change the background to include a cloudy skyline with faint signs of rain approaching.



PLAN1

```
'steps': [{'step_number': 1, 'step_instruction': 'Perspective the train so it appears as if taken from a very close frontal position, emphasizing the front details and giving a sense of proximity to the viewer.',
'instruction_type': 'Reasoning', 'model_choice': 'Bagel'}, {'step_number': 2, 'step_instruction': 'Add a vibrant graffiti art with multicolored abstract designs on the side of the green train's body. Alter the background to depict a cloudy sky with gray, diffused clouds and subtle raindrops starting to fall, giving an ominous yet atmospheric scene.',
'instruction_type': 'Direct',
'model_choice': 'Kontext'}]
```

STEP1



STEP2



Reflect

"quality score": 5

"reflection": "Success" ✅

Figure 2. Visualization of our system's planning, generation and reflection. The model executes the planned editing instructions sequentially. Reflection feedback determines the next action, ensuring better alignment with complex user instructions.

Instruction: Draw what the garden looked like two months ago, replace the watering can with a basket filled with freshly harvested vegetables, and change the sky to a sunset setting.



PLAN1

```
'steps': [{'step_number': 1, 'step_instruction': 'Transition the garden to a state two months prior, showing signs of overgrown plants and faded blooms, with more mature vegetation and slightly wilted plants due to reduced care', 'instruction_type': 'Reasoning', 'model_choice': 'Bagel'}, {'step_number': 2, 'step_instruction': 'replace the watering can with a woven basket filled with freshly harvested vegetables including carrots, lettuce, and tomatoes, and change the sky to a warm sunset hue with vibrant orange and pink clouds', 'instruction_type': 'Direct Alterations', 'model_choice': 'Kontext'}]
```

STEP1



STEP2



Reflect

"quality score": 5  
"reflection": "Success" ✓

Instruction: Replace the background with a mountainous environment. Change the rider's helmet to bright red. Add motion blur to the skateboard and create a trail of dust particles behind it. Introduce another similarly equipped rider slightly behind the main rider.



PLAN1

```
'steps': [{'step_number': 1, 'step_instruction': 'Add another rider similar in appearance and equipment (helmet, attire, skateboard) to the existing rider in the image, positioned slightly behind the main rider to convey a sense of pursuit or companionship.', 'instruction_type': 'Reasoning', 'model_choice': 'Bagel'}, {'step_number': 2, 'step_instruction': 'Replace the background with a photorealistic mountainous environment. Change the current rider's helmet to a bright, saturated red color. Add realistic motion blur to the skateboard to simulate movement, and create a visible trail of dust particles behind the skateboard to emphasize speed and terrain.', 'instruction_type': 'Compositional', 'model_choice': 'Kontext'}]
```

STEP1



STEP2



Reflect

"reason": "The main objectives were mostly achieved: the background was successfully replaced with a mountainous environment, the rider's helmet changed to bright red, motion blur was added to the skateboard, and dust particles were introduced. However, the secondary rider requested to appear slightly behind the main rider is completely missing, which is essential to the instruction."

"quality score": 5  
"reflection": "Refinement" ✗

```
'steps': [{'step_number': 1, 'step_instruction': 'Add a similarly equipped rider—wearing comparable attire and gear—slightly behind the main rider.', 'instruction_type': 'Compositional', 'model_choice': 'Kontext'}]
```

STEP3



Reflect

"quality score": 5  
"reflection": "Success" ✓

Figure 3. Visualization of our system's planning, generation and reflection. The model executes the planned editing instructions sequentially. Reflection feedback determines the next action, ensuring better alignment with complex user instructions.

Instruction	Input	Ours	Bagel-CoT	Bagel	Kontext-Dev	IEAP	ICEdit	OmniGen2	StepDX-Edit	Ovis-U1	CCA
Draw what it will look like after the knot is untied.											
Draw what it will look like after being washed a hundred times.											
Draw what it will look like after being exposed to the sun for one hour.											
Draw what it will look like after three years.											
Generate an image that stacks these blocks in layers from bottom to top, arranged from largest at the base to smallest at the top.											
Generate an image of a fully assembled small table using the provided components with the circular as base and the rectangular as table top.											
Draw what she will look like after a year without cutting her hair.											
Draw what it looks like after being exposed to the sun for a day.											
Draw what it will look like one hour later.											

Figure 4. More visualization results on RISE Bench.

Instruction	Input	Ours	BegeI-CoT	BegeI	Kontext-Dev	IEAP	ICeDit	OmniGen2	StepIX-Edit	Ovis-U1	CCA
Change the fire hydrant's color to bright red. Add a light glow around the fire hydrant. Replace the black metal trash can with a rustic wooden barrel, and surround the base of the fire hydrant with a cluster of flowers. Remove the "1963" inscription in the background.											
Remove the person on the right of the image. Add a duplicate of one of the wine bottles on the counter. Introduce a lit candle on the counter and change the outdoor background to a brightly lit cityscape. Apply a warm-toned filter to harmonize the image.											
Remove the police emblem from the motorcycle. Replace the street background with a beachside boardwalk and add a lens blur while keeping the motorcycle in focus. Change the motorcycle's colors to red and black, and mount a surfboard on it to complete the scene.											
Replace the suburban background with a beach scene, change the frisbee to a larger balloon, change the child's shirt to light blue, and add a flying seagull in the upper right.											
Add a decorative vase with flowers to the table in the background. Remove the glasses from the bed and neatly fold the striped fabric beside where they were. Place an extra pillow on the right side of the bed. Brighten the overall lighting slightly to enhance clarity.											
Replace the backyard background with a field under a clear blue sky, and add a flight of birds in the sky. Change the carrot the subject holds to a bouquet of white and purple flowers, and replace the drink with a cup of coffee.											
Transform the refrigerator's color to pastel pink and adorn its surface with sparkles, while placing a vase with fresh flowers on top of it. Update the wall behind the refrigerator to have a floral-patterned wallpaper, and remove the appliance located to its right.											
Replace the background wall with an artistic mural and change the bunk bed frames to white. Add a hammock hanging diagonally above the bunk beds and modify the bedding texture to feature floral patterns. Adjust the room lighting to a warm golden hue to enhance the atmosphere.											

Figure 5. More visualization results on Complex-Edit Bench.

Instruction	Input	Ours(GPT-4o)	Ours(InternVL)	Nano Banana	GPT-Image	SeedEdt3.0
Draw what it will look like after the <b>car moves out from the parking slot</b> and <b>add a 'Reserved' sign</b> to the parking slot where the car was previously parked. Additionally, <b>change the color of the background</b> to a lighter shade of gray.						
Draw what they will look like when <b>fully inflated</b> , and <b>add a string to each balloon</b> to make them look ready for a party while placing them against a <b>cloudy blue sky background</b> .						
Draw <b>what it looked like during the Jurassic period</b> , and <b>add a lush prehistoric jungle as the background</b> with tall cycads and ferns. Include <b>a flying Pterosaur</b> gliding through the sky above.						reject
Draw an image showing the <b>front view</b> of the house as seen <b>when standing very closely in front of the door</b> . Add a welcome mat with the text <b>'Home Sweet Home'</b> in front of the door and add <b>a flower pot with colorful flowers</b> next to the <b>right side of the door</b> .						
Draw what the mural will look like <b>after six months</b> , accounting for weathering and changes over time, and also <b>add some graffiti text reading 'Dream Big'</b> to the lower-left corner, and reposition <b>the drain pipe to run vertically</b> along the left side of the wall.						
Draw the result of <b>prolonged sunlight exposure</b> on the red chair, and additionally, <b>change the color of the chair to blue</b> and replace the background with an <b>outdoor garden</b> setting.						
Draw what she will look like <b>after a year without cutting her hair</b> , <b>add a necklace</b> with a small pendant around her neck, and place her in a <b>garden setting with vibrant flowers</b> in the background.					reject	
Draw what it will look like when <b>photosynthesis stops for a long time</b> , place the plant on a <b>cracked wooden table</b> , and change the background to a <b>gloomy, overcast outdoor</b> setting.						
Draw what it will look like <b>after a week in a damp basement</b> and <b>add scattered crumbs</b> around the plate along with <b>a faint watermark</b> on the bread that says <b>'For Display Only.'</b>						

Figure 6. More qualitative comparison results of our two variants against several advanced closed-source models on RISE-MIX Bench.