

# INEDIT-BENCH: BENCHMARKING INTERMEDIATE LOGICAL PATHWAYS FOR INTELLIGENT IMAGE EDITING MODELS

**Anonymous authors**

Paper under double-blind review

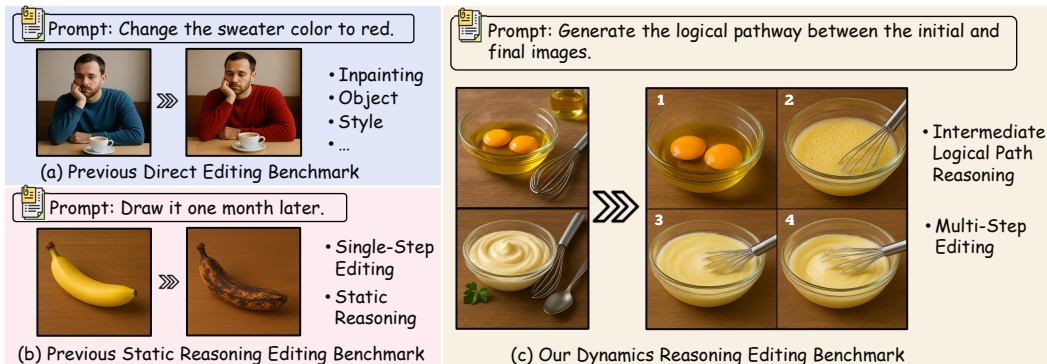


Figure 1: Comparison of previous image editing benchmarks and our proposed InEdit-Bench.

## ABSTRACT

Multimodal generative models have made significant strides in image editing, demonstrating impressive performance on a variety of static tasks. However, their proficiency typically does not extend to complex scenarios requiring dynamic reasoning, leaving them ill-equipped to model the coherent, intermediate logical pathways that constitute a multi-step evolution from an initial state to a final one. This capacity is crucial for unlocking a deeper level of procedural and causal understanding in visual manipulation. To systematically measure this critical limitation, we introduce **InEdit-Bench**, the first evaluation benchmark dedicated to reasoning over intermediate pathways in image editing. InEdit-Bench comprises a meticulously hand-annotated dataset spanning 4 fundamental categories: state transition, dynamic process, temporal sequence, and scientific simulation, which collectively cover 16 distinct sub-tasks. We also propose a suite of 6 evaluation metrics to assess the logical coherence and visual naturalness of the generated pathways, as well as model fidelity to specified or novel path constraints. Our comprehensive evaluation of 14 representative image editing models on InEdit-Bench reveals significant and widespread shortcomings in this domain. By providing a standardized and challenging benchmark, we aim for InEdit-Bench to catalyze research and steer development towards more dynamic, reason-aware, and intelligent multimodal generative models.

## 1 INTRODUCTION

Navigating a complex task is not a single, straightforward jump from inception to conclusion. Instead, the path to the solution is comprised of a series of indispensable intermediate steps that bridge the chasm between the start and the end. Often, the true challenge lies not in the crossing itself, but in the fact that this “bridge” is not readily apparent. We can perceive the starting point and the final destination, yet the pathway connecting them remains invisible. Therefore, the ability to reconstruct this hidden path is a fundamental test of reasoning, prevalent across countless domains.

054 The capacity to reason about transformative pathways is of paramount importance in artificial intel-  
055 ligence, where generative models have unlocked unprecedented prowess in image creation Huang  
056 et al. (2024a); Pan et al. (2025a); Deng et al. (2025). Moreover, intelligent image editing, which  
057 moves beyond simple generation from scratch Fang et al. (2025); Deng et al. (2025), demands a  
058 more profound semantic understanding Wang et al. (2024); Yu et al. (2024) and precise manipula-  
059 tion Alaluf et al. (2021); Brooks et al. (2022), making it a crucial testbed for model competence.  
060 However, despite the remarkable achievements of leading models, they primarily focus on single-  
061 step editing and static reasoning Wu et al. (2025c); Sun et al. (2023); Deng et al. (2025), inherently  
062 lacking the ability to model process evolution. This raises a critical question at the model level:  
063 given only the starting and ending images, how can a model generate a sequence of intermediate  
064 images that adheres to causal logic and visual naturalness?

065 To bring greater attention to this unexplored frontier, we introduce a novel evaluation benchmark,  
066 termed **InEdit-Bench**, centered on the generation of these intermediate logical pathways. As shown  
067 in Fig. 1, our paradigm moves beyond simply appraising the final output. Instead, it challenges a  
068 model to construct the entire, coherent sequence of transformations that logically connects a given  
069 initial state to a final target. This marks a significant departure from existing benchmarks Zhao et al.  
070 (2025); Pan et al. (2025b); Wu et al. (2025d), which, while valuable for assessing static outcomes  
071 like instruction-following fidelity and semantic consistency Pan et al. (2025b); Zhang et al. (2023b),  
072 offer no quantitative measure of procedural reasoning. By shifting the evaluation focus from the  
073 “destination” to the “intermediate logical pathways”, InEdit-Bench provides a more nuanced and  
074 rigorous assessment of the core reasoning faculties, such as causal understanding and strategic plan-  
075 ning. Ultimately, our goal is to steer the research focus away from static, single-step outcomes and  
076 towards the development of models capable of true procedural and dynamic reasoning.

077 Our InEdit-Bench consists of 237 high-quality, meticulously hand-annotated data instances. The  
078 dataset is organized into four fundamental categories: **state transition**, **dynamic process**, **temporal**  
079 **sequence**, and **scientific simulation**, collectively covering 16 distinct sub-tasks. Each instance in  
080 the benchmark comprises an initial state image, a final state image, and a corresponding textual  
081 prompt. To ensure a structured output, these prompts instruct models to generate a single image  
082 divided into  $N$  grids, where each grid depicts a distinct stage of the process. Furthermore, each  
083 prompt contains both a basic editing instruction and a concise summary of key intermediate stages,  
084 generated via a large multimodal model (LMM).

085 For a comprehensive evaluation, InEdit-Bench employs six key metrics to assess the generated pro-  
086 cedural pathways: **appearance consistency**, **perceptual quality**, **semantic consistency**, **logical co-**  
087 **herence**, **scientific plausibility**, and **process plausibility**. While the first three metrics are adapted  
088 from standard image editing tasks Pan et al. (2025b); Zhao et al. (2025), the latter three are novel  
089 and specifically designed for our process-oriented benchmark. These new metrics provide an ob-  
090 jective assessment of the transition logic between stages, their scientific fidelity, and the holistic  
091 comprehension of the intermediate pathway. To automate this multifaceted evaluation, we adopt  
092 the LMM-as-a-Judge paradigm Zhao et al. (2025), where a powerful LMM serves as an objective  
093 evaluator for the generated image pathways.

094 Using InEdit-Bench, we conduct an evaluation of 14 representative intelligent image editing meth-  
095 ods, including state-of-the-art models such as GPT-Image-1 OpenAI (2025), Nano-Banana Google  
096 (2025), Flux-Kontext-pro Labs (2025), Qwen-Image-Edit Wu et al. (2025a), Bagel Deng et al.  
097 (2025), Emu Sun et al. (2023), OmniGen Wu et al. (2025c), and Step1X-Edit Liu et al. (2025).  
098 Our findings consistently reveal that these models exhibit significant shortcomings in multi-step  
099 editing and dynamic reasoning, highlighting a crucial direction for future research.

100 In summary, our main contributions are as follows:

101 (1) We introduce InEdit-Bench, the first evaluation benchmark for multi-step image editing and  
102 dynamic reasoning. It provides a challenging testbed to assess the ability of a model to comprehend  
103 and generate intermediate logical pathways, catalyzing future research in controllable visual editing.

104 (2) We construct a high-quality, meticulously annotated dataset for the benchmark, encompassing  
105 4 fundamental categories and 16 distinct sub-tasks. This dataset establishes a robust foundation for  
106 the systematic evaluation of complex editing capabilities.

107

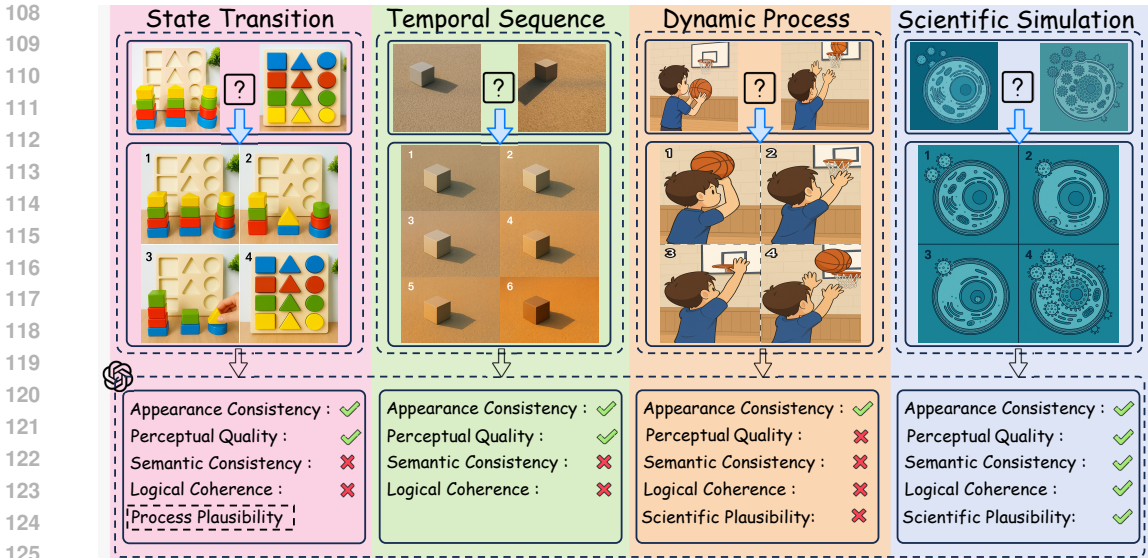


Figure 2: **Overall introduction to InEdit-Bench.** InEdit-Bench focuses on dynamic reasoning and multi-step editing modes, requiring models to generate intermediate logical pathways for given tasks. It spans 4 key domains: state transition, dynamic process, temporal sequence, and scientific simulation. The evaluation is conducted through 6 dimensions: appearance consistency, perceptual quality, semantic consistency, logical coherence, scientific plausibility, and process plausibility.

(3) We establish a multi-faceted evaluation protocol with 6 assessment dimensions. These metrics are specifically designed to capture the visual fidelity and logical coherence of intermediate paths, ensuring a rigorous and objective assessment.

(4) We present a comprehensive analysis of 14 state-of-the-art models on InEdit-Bench. Our findings reveal the significant limitations of current methods and highlight key areas for future improvement.

## 2 RELATED WORK

### 2.1 INSTRUCTION-BASED IMAGE EDITING

Instruction-based Image Editing (IIE) Wang et al. (2024); Huang et al. (2024a); Pan et al. (2025a) emphasizes the direct expression of user intent, thereby reducing interaction complexity while improving controllability and practicality. Early work such as InstructPix2Pix Brooks et al. (2022) proposed driving image editing with simple instructions, and subsequent research has continuously improved data quality and model architectures. For example, MagicBrush Zhang et al. (2023a) introduced high-quality manually annotated data, MGIE Fu et al. (2024) and SmartEdit Huang et al. (2024b) leveraged multimodal large language models to enhance semantic understanding and reasoning, while OmniGen Wu et al. (2025c) and Gemini Team & et al. (2025) built unified multimodal architectures to strengthen task generalization. Although existing methods have achieved progress in quality Podell et al. (2024); Rombach et al. (2022), efficiency Team (2024); Wu et al. (2025b), and controllability Du et al. (2025), the capabilities of models in multi-step editing and higher-order understanding remain underexplored. To address this, we propose InEdit-Bench, aiming to provide valuable insights for the community and advance the development of such models.

### 2.2 IMAGE EDITING BENCHMARKS

The proliferation of image editing benchmarks has accelerated the development of model evaluation frameworks in recent years. Benchmarks such as TedBench Kawar et al. (2023), EditVal Basu et al. (2023), and EditBench Wang et al. (2023) primarily focus on basic editing tasks, while MagicBrush Zhang et al. (2023a) provides high-quality data but its evaluation metrics Caron et al. (2021); Radford et al. (2021) have inherent limitations in fully reflecting image quality. With further research, Reason-edit Huang et al. (2024b) and RISEBench Zhao et al. (2025) have begun to em-

phasize evaluating models’ understanding and reasoning capabilities under complex instructions. Complex-Edit Yang et al. (2025) introduces a chain-of-thought-like multi-step editing mechanism, while I2EBench Ma et al. (2024) and KRIS-Bench Wu et al. (2025d) explore higher-level editing abilities from the perspectives of multi-dimensional skills and knowledge-driven reasoning. Overall, although existing benchmarks have made significant progress in terms of task scale and diversity, they remain insufficient in modeling intermediate logical pathways and supporting dynamic chain-style reasoning. To address this gap, we propose InEdit-Bench, a new benchmark for intermediate logical pathway editing, designed to systematically evaluate the capabilities of intelligent visual editing models in multi-step image editing as well as dynamic understanding and reasoning.

### 3 INEDIT-BENCH

#### 3.1 BENCHMARK CONSTRUCTION

InEdit-Bench is an innovative benchmark designed to systematically evaluate a model’s capability in comprehending and representing intermediate logical pathways. As illustrated in Fig. 2, we categorize editing tasks into four fundamental types based on the evolutionary dynamics of their intermediate states: **state transition**, **dynamic process**, **temporal sequence**, and **scientific simulation**. Fig. 3 details the task distribution within InEdit-Bench. The representative example images of each subtask are provided in Appx. A.6. The image sources for this benchmark include data collected from the internet under permissive licenses, images generated by generative models OpenAI (2025), and samples extracted from existing datasets Wu et al. (2025d); Lanitis et al. (2002).

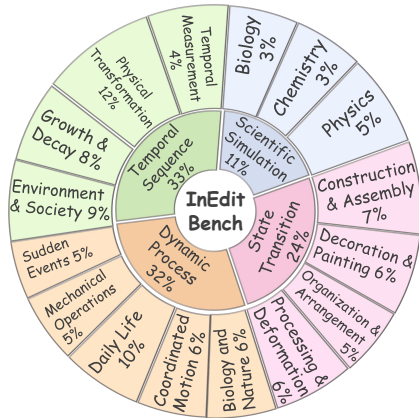


Figure 3: The task type distribution of InEdit-Bench. InEdit-Bench conducts a comprehensive evaluation of visual editing models across 16 sub-tasks under 4 domains.

##### 3.1.1 STATE TRANSITION

State transition reasoning aims to understand and reconstruct the discrete changes from an initial state to a final state. The primary challenge is to identify key discrete change nodes and infer their dependencies, thereby generating a complete editing sequence that is both logically coherent and structurally clear. This task category is further subdivided into the following four subcategories. (1) **Construction and Assembly:** This subtask requires the model to combine multiple independent components into a complete entity based on spatial logic. The main difficulty lies in structural reconstruction and resolving the dependencies between assembly steps. (2) **Decoration and Painting:** This pertains to applying colors, textures, or patterns to a target’s surface. The critical aspect is achieving precise region identification and control over the sequence of operations. (3) **Organization and Arrangement:** The central difficulty here is the systematic organization and spatial arrangement of elements, which requires the model to possess strong reasoning capabilities regarding layout structures and spatial relationships. (4) **Processing and Deformation:** This involves the model understanding how an operation alters the physical properties or state of an object. The key challenge is to precisely capture and simulate the object’s morphological transformations.

##### 3.1.2 DYNAMIC PROCESS

Dynamic process reasoning is characterized by a continuous transformation from an initial state to a final one. In contrast to scenarios involving discrete state transition, this paradigm challenges the model to process seamless and continuous progressions, demanding that every intermediate step demonstrates both natural fluidity and logical consistency. The tasks within this category are further classified into five distinct subdomains. (1) **Biology and Nature:** Focuses on the evolution of organisms and natural phenomena, stressing deduction from biological characteristics and the laws of nature (e.g., a spider constructing a web, a chick hatching). (2) **Coordinated Motion:** Concerns the fluid and coordinated movement of entities through space, mandating that the model comprehend

216 motion dynamics to produce smooth, logically connected intermediate actions (*e.g.*, a long jump, a  
 217 basketball shot). **(3) Daily Life:** Encompasses the modeling of common operations, interactions,  
 218 and behaviors observed in everyday contexts. **(4) Mechanical Operation:** Involves illustrating the  
 219 continuous structural alteration of an object via incremental reasoning, highlighting the operational  
 220 mechanisms that propel such changes (*e.g.*, a compressor flattening a cube). **(5) Sudden Events:**  
 221 Characterized by abrupt, often destructive transformations, requiring the model to identify pivotal  
 222 moments and render believable visual consequences (*e.g.*, a building demolition). To systemati-  
 223 cally guide the model in generating these continuous transformations, the task instructions for each  
 224 instance provide one to three key path-node prompts.

### 225 226 3.1.3 TEMPORAL SEQUENCE

228 Temporal sequence reasoning is concerned with the progressive evolution of a target state over time.  
 229 Diverging from the emphasis on continuity in dynamic process, temporal tasks prioritize identify-  
 230 ing and demarcating critical change points on a timeline. Such tasks necessitate that the model is  
 231 endowed with a capacity for temporal-aware inference and can accurately represent the distinct char-  
 232 acteristics of each evolutionary phase. The editing protocol requires the model to partition the entire  
 233 process into uniform temporal intervals, thus providing a clear and complete trajectory of the tem-  
 234 poral evolution. This domain is subdivided into the following four subcategories. **(1) Environment  
 235 and Society:** Concerns the progressive evolution of environmental phenomena and social behaviors  
 236 (*e.g.*, a train arriving at a station, the formation of sand ripples), requiring the model to reason about  
 237 and generate temporally congruent sequences for such events. **(2) Growth and Decay:** Pertains  
 238 to the life cycles of organisms (*e.g.*, a flower blooming, a wound healing), requiring the generation  
 239 of a time-series that aligns with biological principles. **(3) Physical Transformation:** Typically in-  
 240 volves alterations in the physical properties of materials or objects (*e.g.*, ice melting). The primary  
 241 difficulty lies in modeling the temporal progression of these properties while ensuring the logical  
 242 coherence of the transitional states. **(4) Temporal Measurement:** Focuses on the representation of  
 243 time as a quantifiable metric (*e.g.*, a progress bar, an hourglass), demanding precise reasoning about  
 244 quantitative changes along the temporal dimension.

### 245 246 3.1.4 SCIENTIFIC SIMULATION

248 Scientific simulation is designed to model principles from the fields of **Physics, Chemistry, and  
 249 Biology**. These domains impose a strict requirement on the model to adhere to scientific laws, while  
 250 illuminating the intermediate logical steps of a given process. For instance, physical phenomena  
 251 range from the diffusion effect to total internal reflection; chemical processes include reactions such  
 252 as the combustion of a magnesium strip and displacement reactions; and biological principles are  
 253 exemplified by life processes like cell division and DNA replication. These tasks require the model  
 254 to comprehend and execute complex scientific procedures, deduce causal mechanisms, and render  
 255 each pivotal stage. To streamline this process, task instructions provide concise keyframe prompts,  
 256 ensuring the model concentrates on the most critical phases and disregards superfluous steps.

## 257 258 3.2 EVALUATION METRICS

260 Diverging from conventional benchmarks that assess single-step image editing, InEdit-Bench funda-  
 261 mentally redefines the evaluation paradigm by reorienting the focus from input-output comparisons  
 262 to the procedural integrity of the entire transformation. To this end, we develop a multi-faceted  
 263 evaluation framework built upon six key dimensions. These are categorized into two groups: three  
 264 foundational metrics for visual quality (**Appearance Consistency, Perceptual Quality, Semantic  
 265 Consistency**), and three novel dimensions (see Fig. 4) designed to scrutinize the plausibility of in-  
 266 termediate processes (**Logical Coherence, Scientific Plausibility, Process Plausibility**). Our eval-  
 267 uation employs the LMM-as-a-Judge methodology, leveraging GPT-4o for its advanced reasoning  
 268 capabilities, which are essential for our novel process-oriented metrics. For evaluation, the model  
 269 receives the user instruction, a scoring rubric, and the generated output—a single image depicting  
 the process across a grid, based on which it provides a numerical score for each dimension.

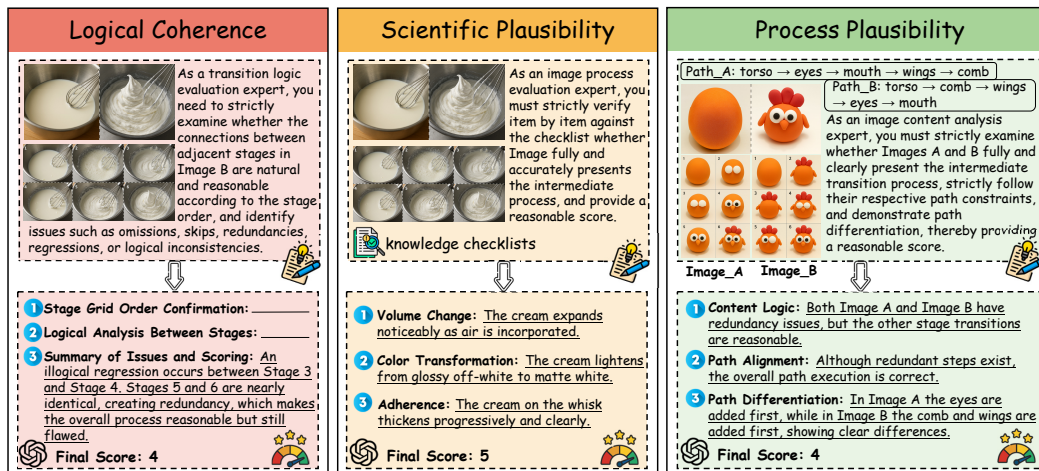


Figure 4: The evaluation metrics of **Logical Coherence**, **Scientific Plausibility**, and **Process Plausibility** in InEdit-Bench.

### 3.3 STANDARD VISUAL QUALITY METRICS

To establish a baseline for visual quality, our framework incorporates three foundational metrics from the image editing domain. Appearance consistency assesses the preservation of style and visual attributes across all depicted stages of the process. Perceptual quality measures the realism and fidelity of the generated imagery, ensuring it is free from artifacts. Semantic consistency evaluates the alignment of the final image content with the specified editing objective. Collectively, these metrics provide a robust assessment of the visual integrity of the final output.

A critical guideline for our benchmark stems from the use of a grid-based representation. The evaluation protocol explicitly requires that visual artifacts introduced by the grid format (such as grid lines, segmentation effects, and numbering) must be disregarded during content assessment. Furthermore, any layout or compositional discrepancies arising as a direct consequence of the grid structure are also to be excluded from the evaluation.

### 3.4 PROPOSED PROCESS-ORIENTED METRICS

**Logical Coherence.** Logical coherence is paramount in evaluating multi-step image editing, as it examines the integrity of the generated process in terms of its logical progression and natural flow. The assessment protocol begins by establishing the sequence of the depicted stages, applying a top-to-bottom, left-to-right convention whenever the intended order is not visually apparent. The core of the evaluation then involves a close examination of the transitions between adjacent stages for logical soundness and naturalness. This scrutiny ensures that the overall evolution is fluid and coherent, devoid of any jarring discontinuities or superfluous, repetitive actions.

**Scientific Plausibility.** Drawing inspiration from KRIS-Bench Wu et al. (2025d) and WorldGen-Bench Zhang et al. (2025), we incorporate scientific plausibility as a dedicated evaluation dimension. This metric is applied to tasks involving dynamic process and scientific simulation, where adherence to scientific logic is assessed via a series of knowledge checklists. These checklists meticulously annotate the critical features and inherent mechanisms that should be present in the intermediate stages. The assessment conducts a direct comparison of the generated visual content against the checklist’s items to verify compliance with the predefined scientific standards.

**Process Plausibility.** For a comprehensive evaluation of how well a model understands intermediate pathways, we employ two prompting schemes with distinct path constraints, applied to a subset of our state transition tasks. This approach is motivated by the fact that many real-world processes are non-deterministic, often presenting multiple viable routes to the same outcome. Consequently, a capable model must not only grasp the fundamental operation of each step but also discern a rational sequence from among multiple viable paths, all while maintaining both consistency and accuracy.

Table 1: **Performance of various models on InEdit-Bench.** For the dynamic process and scientific simulation tasks, scores in gray denote performance calculated without the *scientific plausibility* metric. For the state transition task, scores in gray denote the model’s performance on *process plausibility* data. The best and second-best results are highlighted in **bold** and with an underline, respectively.

Models	State Transition	Temporal Sequence	Dynamic Process	Scientific Simulation	Overall Average	Accuracy
<b>Proprietary Models</b>						
GPT-Image-1	<b>81.12</b> [89.00]	<b>81.25</b>	<b>79.85</b> (82.21)	<b>82.61</b> (84.78)	<b>81.33</b>	<b>16.75%</b>
Nano-Banana	70.15 [75.70]	74.24	77.85 (79.33)	79.57 (79.89)	75.23	13.30%
Flux-Kontext-pro	51.76 [47.10]	52.18	56.23 (59.33)	51.30 (53.26)	51.46	0.99%
Doubao-SeedEdit-3.0	38.36 [25.00]	42.23	39.54 (39.81)	33.48 (35.60)	36.50	0.00%
<b>Open-Source Models</b>						
Qwen-Image-Edit	<b>44.77</b> [51.50]	<b>52.08</b>	<b>52.92</b> (54.23)	<b>44.13</b> (45.92)	<b>49.60</b>	<u>0.49%</u>
Emu1	12.63 [3.70]	14.96	14.54 (17.02)	12.39 (13.86)	11.42	0.00%
Emu2	33.46 [15.40]	34.55	34.08 (37.31)	31.30 (32.61)	29.61	0.00%
Bagel	38.27 [42.60]	42.61	42.31 (44.33)	39.13 (40.49)	<u>40.70</u>	0.00%
Bagel-Think	<u>41.84</u> [26.50]	<u>44.79</u>	<u>46.92</u> (49.81)	<u>43.48</u> (46.74)	<u>40.70</u>	<b>0.99%</b>
OmniGen	10.46 [14.00]	16.29	16.15 (16.54)	12.39 (12.50)	14.34	0.00%
OmniGen2	38.39 [30.90]	43.47	41.85 (44.42)	34.35 (37.23)	37.92	<u>0.49%</u>
Step1X-Edit(v1.0)	17.18 [9.60]	17.23	20.99 (22.11)	13.91 (14.67)	16.43	0.00%
Step1X-Edit(v1.1)	23.47 [26.50]	32.20	37.54 (39.13)	33.91 (33.97)	31.39	0.00%
InstructPix2Pix	28.57 [0.00]	37.07	26.38 (29.33)	24.35 (27.99)	23.23	0.00%

Accordingly, we introduce process plausibility as an advanced metric that evaluates a model’s comprehension of intermediate logical pathways. Under this metric, a successful generation must adhere to the global path sequence and avoid internal logical errors or representational deviations, thereby ensuring a clear and accurate depiction of the intended logical trajectory. A comparative analysis of model performance under the two constraint schemes reveals the overall grasp of path constraints and the capability to follow a specified path.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTS SETUP

On InEdit-Bench, we evaluated 14 representative visual editing models, conducting a detailed quantitative analysis of their performance across six dimensions. The evaluated models include proprietary models: GPT-Image-1 OpenAI (2025), Nano-Banana Google (2025), Flux-Kontext-pro Labs (2025), and Doubao-SeedEdit-3.0-i2i ByteDance (2025); as well as open-source models: Qwen-Image-Edit Wu et al. (2025a), Bagel Deng et al. (2025), OmniGen Xiao et al. (2025), OmniGen2 Wu et al. (2025c), Step1X-Edit Liu et al. (2025), Emu1 Sun et al. (2024), Emu2 Sun et al. (2023), and InstructPix2Pix Brooks et al. (2022). These models cover a range of mainstream generative architectures, including autoregressive generation paradigms Sun et al. (2023); Deng et al. (2025), diffusion model architectures Brooks et al. (2022), and diffusion transformer architectures Liu et al. (2025). All generation and evaluation processes are conducted on L20 GPUs using default hyperparameter settings to ensure fairness and reproducibility. Additionally, since some open-source models do not support multi-image input, we uniformly concatenated the initial and final state images into a single image, with the initial state image placed on the left or top, and the final state image on the right or bottom, separated by a black-and-white striped line.

### 4.2 RESULT ANALYSIS

#### 4.2.1 RESULTS ANALYSIS BY TASKS

Tab. 1 reports the scores of 14 models on four fundamental tasks. All scores are normalized to a 100-point scale and are evaluated by GPT-4o-2024-11-20. Results on InEdit-Bench show that GPT-Image-1 is the best-performing proprietary model, with a score of 81.33 and an accuracy of 16.75%. Among open-source models, Qwen-Image-Edit and Bagel-Think perform relatively well, scoring

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

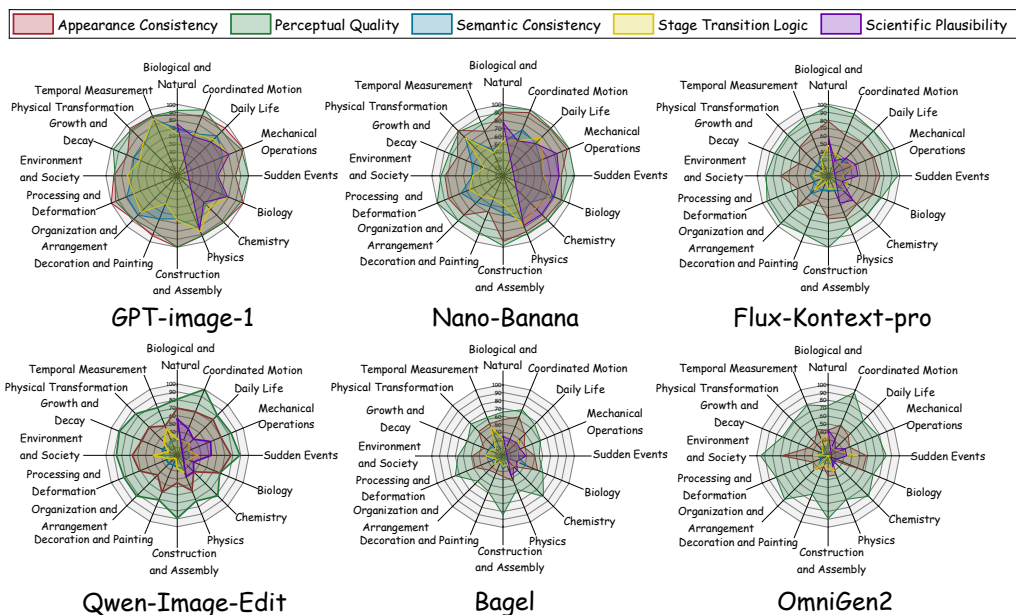


Figure 5: The performance of GPT-Image-1, Nano-Banana, Flux-Kontext-pro, Qwen-Image-Edit, Bagel and OmniGen2 across 16 sub-tasks.

49.60 and 40.70 points, respectively. Although there is still a gap between open-source and proprietary models, some open-source models have nevertheless demonstrated outstanding capabilities.

From the perspective of task dimensions, state transition task poses significant challenges for models. The average scores of all models in this category are lower than their performance on temporal sequence and dynamic process tasks. Notably, nine models, including GPT-Image-1 and Nano-Banana, score lower on state transition task than on scientific simulation task (when comparing state transition, dynamic process, and scientific simulation tasks, only the average scores of four metrics—appearance consistency, perceptual quality, semantic consistency, and logical coherence—are uniformly calculated). Furthermore, apart from the generally stronger GPT-Image-1 and Nano-Banana, all other models achieve lower average scores on scientific simulation task compared to their performance on dynamic process task. These phenomena reveal a layered structure of task complexity: from continuous to discrete, and from surface-level phenomena to deeper scientific principles, model performance shows a step-by-step decline. This highlights the limitations of current large models in complex logical reasoning and scientific law modeling.

In terms of process plausibility, GPT-Image-1 achieves the highest score of 89.00, demonstrating its advantage in understanding and articulating reasoning paths, with Nano-Banana ranking closely behind. In contrast, most of the other models generally struggle to fully meet task requirements. With respect to accuracy, even the best-performing model, GPT-Image-1, attains only 16.75%, followed by Nano-Banana at 13.30%. The remaining models achieve accuracies below 1.00%, with only Flux-Kontext-pro, Qwen-Image-Edit, Bagel-Think, and OmniGen2 reaching 0.99%, 0.49%, 0.99%, and 0.49%, respectively. Most other models yield 0% accuracy. Overall, these results reveal that current models still face significant limitations in long-term dependency capture and multi-stage causal reasoning. Achieving accurate modeling and representation of intermediate reasoning paths remains a key challenge that urgently needs to be addressed in the future.

Fig. 5 presents the performance of several representative models on 16 subtasks, while the complete results for the remaining models are provided in Appx. A.2. Overall, GPT-Image-1 demonstrates stable and superior performance across all subtasks. In contrast, the results of Flux-Kontext-pro, Qwen-Image-Edit, Bagel, and OmniGen2 exhibit significant fluctuations, particularly in the dimensions of semantic consistency and logical coherence, where substantial performance gaps may arise even among subtasks within the same basic category. Specifically, within the temporal sequence category, Flux-Kontext-pro and Bagel experience a notable drop in performance on the physical transformation subtask, while Qwen-Image-Edit and OmniGen2 show marked degradation on the

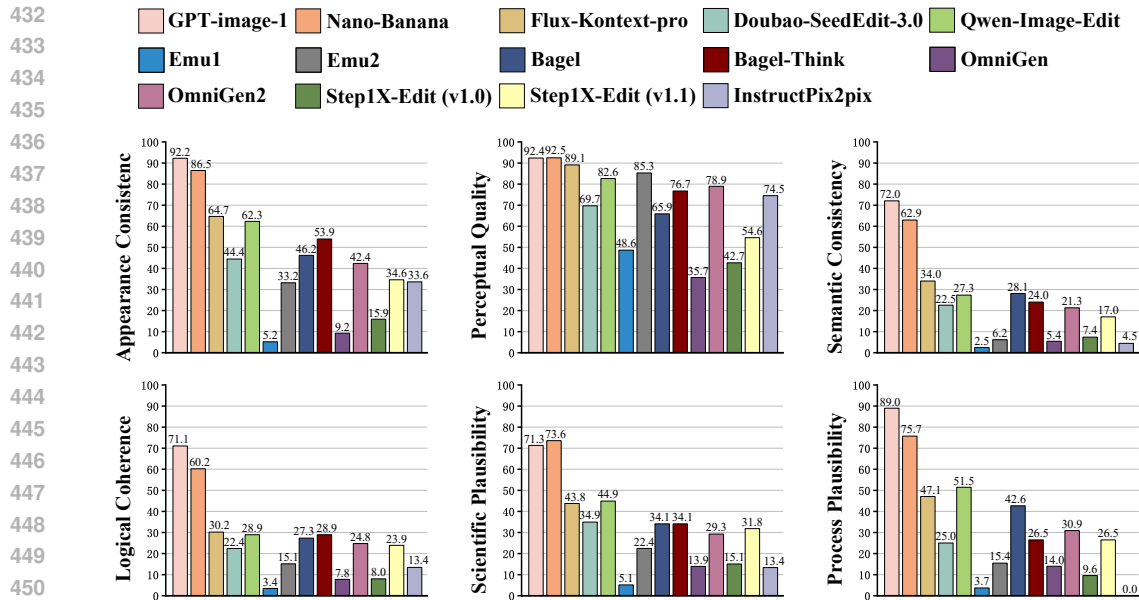


Figure 6: Comparison of models across six evaluation dimensions.

growth and decay subtask. Nano-Banana achieves relatively better average performance on dynamic process and scientific simulation tasks. However, it still suffers from performance decline and uneven results in the state transition and temporal sequence categories.

#### 4.2.2 RESULTS ANALYSIS BY METRICS

As illustrated in Fig. 6, our analysis indicates that proprietary models, specifically GPT-Image-1 and Nano-Banana, consistently lead in appearance consistency, semantic consistency, and logical coherence, demonstrating their robust and well-rounded capabilities. While open-source models post lower aggregate scores, several demonstrate significant potential in specific dimensions. For instance, Qwen-Image-Edit stands out among open-source solutions for the high performance in semantic consistency, logical coherence, and scientific plausibility, occasionally rivaling proprietary counterparts. Similarly, the Bagel series is highly competitive in perceptual quality and semantic consistency, a strength also exhibited by OmniGen2 in perceptual quality.

On the other hand, open-source models Emu1 and OmniGen struggle to maintain effective visual consistency. Among proprietary models, Doubao lags behind, with significantly lower scores in appearance consistency, semantic consistency, and logical coherence compared with its peers. This suggests that Doubao may place more emphasis on rapid local editing while lacking robustness in modeling global consistency and cross-modal logical constraints. Notably, half of the open-source models score below 10.00 in the semantic consistency dimension, further underscoring their systematic deficiencies in intermediate logical path editing, instruction understanding, and semantically effective editing.

## 5 CONCLUSION

In this paper, we propose InEdit-Bench, the first system evaluation benchmark focused on image multi-step editing and intermediate logical path reasoning. It covers four basic categories: state transition, dynamic process, temporal sequence, and scientific simulation, with 16 sub-tasks. Based on this, we design six evaluation dimensions: appearance consistency, perceptual quality, semantic consistency, logical coherence, scientific plausibility, and process plausibility, to comprehensively measure the ability of intelligent visual editing models in intermediate logical path reasoning and expression. Through the evaluation of 14 representative models, we reveal significant shortcomings in current models regarding multi-step editing and dynamic reasoning capabilities, providing clear directions and reference for further optimization of model performance.

## REFERENCES

- 486  
487  
488 Yuval Alaluf, Omer Tov, Ron Mokady, Rinon Gal, and Amit H. Bermano. Hyperstyle: Style-  
489 gan inversion with hypernetworks for real image editing. *2022 IEEE/CVF Conference on*  
490 *Computer Vision and Pattern Recognition (CVPR)*, pp. 18490–18500, 2021. URL <https://api.semanticscholar.org/CorpusID:244729249>.  
491
- 492 Samyadeep Basu, Mehrdad Saberi, Shweta Bhardwaj, Atoosa Malemir Chegini, Daniela Massiceti,  
493 Maziar Sanjabi, Shell Xu Hu, and Soheil Feizi. Editval: Benchmarking diffusion based text-  
494 guided image editing methods. *CoRR*, abs/2310.02426, 2023. doi: 10.48550/ARXIV.2310.02426.  
495 URL <https://doi.org/10.48550/arXiv.2310.02426>.  
496
- 497 Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow im-  
498 age editing instructions. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recog-  
499 nition (CVPR)*, pp. 18392–18402, 2022. URL [https://api.semanticscholar.org/  
500 CorpusID:253581213](https://api.semanticscholar.org/CorpusID:253581213).
- 501 ByteDance. Doubao-seededit-3.0-i2i-250628. <https://console.volcengine.com/ark>, 2025. 2025-09-  
502 01.
- 503 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and  
504 Armand Joulin. Emerging properties in self-supervised vision transformers. In *2021 IEEE/CVF  
505 International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-  
506 17, 2021*, pp. 9630–9640. IEEE, 2021. doi: 10.1109/ICCV48922.2021.00951. URL <https://doi.org/10.1109/ICCV48922.2021.00951>.  
507  
508
- 509 Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Wei-  
510 hao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified  
511 multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025.
- 512 Yifan Du, Hangyu Guo, Kun Zhou, Wayne Xin Zhao, Jinpeng Wang, Chuyuan Wang, Mingchen  
513 Cai, Ruihua Song, and Ji-Rong Wen. What makes for good visual instructions? synthesizing  
514 complex visual reasoning instructions for visual instruction tuning. In Owen Rambow, Leo Wan-  
515 ner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert (eds.),  
516 *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025,  
517 Abu Dhabi, UAE, January 19-24, 2025*, pp. 8197–8214. Association for Computational Linguis-  
518 tics, 2025. URL <https://aclanthology.org/2025.coling-main.546/>.
- 519 Rongyao Fang, Chengqi Duan, Kun Wang, Linjiang Huang, Hao Li, Shilin Yan, Hao Tian, Xingyu  
520 Zeng, Rui Zhao, Jifeng Dai, Xihui Liu, and Hongsheng Li. Got: Unleashing reasoning ca-  
521 pability of multimodal large language model for visual generation and editing. *arXiv preprint  
522 arXiv:2503.10639*, 2025.  
523
- 524 Tsu-Jui Fu, Wenze Hu, Xianzhi Du, William Yang Wang, Yinfei Yang, and Zhe Gan. Guiding  
525 instruction-based image editing via multimodal large language models. In *The Twelfth Interna-  
526 tional Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.  
527 OpenReview.net, 2024. URL <https://openreview.net/forum?id=S1RKWSyZ2Y>.
- 528 Google. Gemini 2.5 flash image preview. <https://ai.google.dev/gemini-api/docs/models>, 2025. 2025-  
529 09-01.  
530
- 531 Yi Huang, Jiancheng Huang, Yifan Liu, Mingfu Yan, Jiaxi Lv, Jianzhuang Liu, Wei Xiong,  
532 He Zhang, Shifeng Chen, and Liangliang Cao. Diffusion model-based image editing: A sur-  
533 vey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47:4409–4437, 2024a.  
534 URL <https://api.semanticscholar.org/CorpusID:268033671>.
- 535 Yuzhou Huang, Liangbin Xie, Xintao Wang, Ziyang Yuan, Xiaodong Cun, Yixiao Ge, Jiantao  
536 Zhou, Chao Dong, Rui Huang, Ruimao Zhang, and Ying Shan. Smartedit: Exploring com-  
537 plex instruction-based image editing with multimodal large language models. In *IEEE/CVF  
538 Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June  
539 16-22, 2024*, pp. 8362–8371. IEEE, 2024b. doi: 10.1109/CVPR52733.2024.00799. URL  
<https://doi.org/10.1109/CVPR52733.2024.00799>.

- 540 Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri,  
541 and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *IEEE/CVF*  
542 *Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada,*  
543 *June 17-24, 2023*, pp. 6007–6017. IEEE, 2023. doi: 10.1109/CVPR52729.2023.00582. URL  
544 <https://doi.org/10.1109/CVPR52729.2023.00582>.
- 545 Black Forest Labs. Flux-kontext-pro. [https://fluxproweb.com/cn/blog/detail/Introducing-Flux-Pro-](https://fluxproweb.com/cn/blog/detail/Introducing-Flux-Pro-Kontext2025-09-01)  
546 [Kontext2025-09-01](https://fluxproweb.com/cn/blog/detail/Introducing-Flux-Pro-Kontext2025-09-01).
- 547 Andreas Lanitis, Christopher J. Taylor, and Timothy F Cootes. Toward automatic simulation of  
548 aging effects on face images. *IEEE Transactions on pattern Analysis and machine Intelligence*,  
549 24(4):442–455, 2002.
- 550 Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming  
551 Wang, Honghao Fu, Chunrui Han, Guopeng Li, Yuang Peng, Quan Sun, Jingwei Wu, Yan Cai,  
552 Zheng Ge, Ranchen Ming, Lei Xia, Xianfang Zeng, Yibo Zhu, Binxing Jiao, Xiangyu Zhang,  
553 Gang Yu, and Daxin Jiang. Step1x-edit: A practical framework for general image editing. *arXiv*  
554 *preprint arXiv:2504.17761*, 2025.
- 555 Yiwei Ma, Jiayi Ji, Ke Ye, Weihuang Lin, Zhibin Wang, Yonghan Zheng, Qiang Zhou, Xi-  
556 aoshuai Sun, and Rongrong Ji. I2ebench: A comprehensive benchmark for instruction-  
557 based image editing. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela  
558 Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in Neural*  
559 *Information Processing Systems 38: Annual Conference on Neural Information Process-*  
560 *ing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15,*  
561 *2024*, 2024. URL [http://papers.nips.cc/paper\\_files/paper/2024/hash/](http://papers.nips.cc/paper_files/paper/2024/hash/48fecef47b19fe501d27d338b6d52582-Abstract-Conference.html)  
562 [48fecef47b19fe501d27d338b6d52582-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2024/hash/48fecef47b19fe501d27d338b6d52582-Abstract-Conference.html).
- 563 OpenAI. Gpt-image-1: Openai’s multimodal image generation model.  
564 <https://platform.openai.com/docs/models/gpt-image-1>, 2025. 2025-09-01.
- 565 Kaihang Pan, Wang Lin, Zhongqi Yue, Tenglong Ao, Liyu Jia, Wei Zhao, Juncheng Li, Siliang  
566 Tang, and Hanwang Zhang. Generative multimodal pretraining with discrete diffusion timestep  
567 tokens. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025,*  
568 *Nashville, TN, USA, June 11-15, 2025*, pp. 26136–26146. Computer Vision Foundation / IEEE,  
569 2025a. URL [https://openaccess.thecvf.com/content/CVPR2025/html/](https://openaccess.thecvf.com/content/CVPR2025/html/Pan_Generative_Multimodal_Pretraining_with_Discrete_Diffusion_Timestep_Tokens_CVPR_2025_paper.html)  
570 [Pan\\_Generative\\_Multimodal\\_Pretraining\\_with\\_Discrete\\_Diffusion\\_](https://openaccess.thecvf.com/content/CVPR2025/html/Pan_Generative_Multimodal_Pretraining_with_Discrete_Diffusion_Timestep_Tokens_CVPR_2025_paper.html)  
571 [Timestep\\_Tokens\\_CVPR\\_2025\\_paper.html](https://openaccess.thecvf.com/content/CVPR2025/html/Pan_Generative_Multimodal_Pretraining_with_Discrete_Diffusion_Timestep_Tokens_CVPR_2025_paper.html).
- 572 Yulin Pan, Xiangteng He, Chaojie Mao, Zhen Han, Zeyinzi Jiang, Jingfeng Zhang, and Yu Liu.  
573 Ice-bench: A unified and comprehensive benchmark for image creating and editing. *ArXiv*,  
574 [abs/2503.14482](https://arxiv.org/abs/2503.14482), 2025b. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:277104618)  
575 [277104618](https://api.semanticscholar.org/CorpusID:277104618).
- 576 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
577 Penna, and Robin Rombach. SDXL: improving latent diffusion models for high-resolution image  
578 synthesis. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vi-*  
579 *enna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=di5zR8xgf)  
580 [forum?id=di5zR8xgf](https://openreview.net/forum?id=di5zR8xgf).
- 581 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
582 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
583 Sutskever. Learning transferable visual models from natural language supervision. In Marina  
584 Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Ma-*  
585 *chine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Ma-*  
586 *chine Learning Research*, pp. 8748–8763. PMLR, 2021. URL [http://proceedings.mlr.](http://proceedings.mlr.press/v139/radford21a.html)  
587 [press/v139/radford21a.html](http://proceedings.mlr.press/v139/radford21a.html).
- 588 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
589 resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Com-*  
590 *puter Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*,  
591 pp. 10674–10685. IEEE, 2022. doi: 10.1109/CVPR52688.2022.01042. URL [https://doi.](https://doi.org/10.1109/CVPR52688.2022.01042)  
592 [org/10.1109/CVPR52688.2022.01042](https://doi.org/10.1109/CVPR52688.2022.01042).
- 593

- 594 Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Zhengxiong Luo, Yueze Wang,  
595 Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models  
596 are in-context learners. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14398–14409, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:266374640)  
597 [CorpusID:266374640](https://api.semanticscholar.org/CorpusID:266374640).  
598
- 599 Quan Sun, Qiyang Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao,  
600 Jingjing Liu, Tiejun Huang, and Xinlong Wang. Emu: Generative pretraining in multimodality. In  
601 *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria,*  
602 *May 7-11, 2024*. OpenReview.net, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=mL8Q900amV)  
603 [mL8Q900amV](https://openreview.net/forum?id=mL8Q900amV).  
604
- 605 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *CoRR*,  
606 [abs/2405.09818](https://arxiv.org/abs/2405.09818), 2024. doi: 10.48550/ARXIV.2405.09818. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2405.09818)  
607 [48550/arXiv.2405.09818](https://doi.org/10.48550/arXiv.2405.09818).  
608
- 609 Gemini Team and Rohan Anil et al. Gemini: A family of highly capable multimodal models, 2025.  
610 URL <https://arxiv.org/abs/2312.11805>.
- 611 Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Ya-  
612 sumasa Onoe, Sarah Laszlo, David J. Fleet, Radu Soricut, Jason Baldridge, Mohammad Norouzi,  
613 Peter Anderson, and William Chan. Imagen editor and editbench: Advancing and evaluating  
614 text-guided image inpainting. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pp. 18359–18369. IEEE, 2023.  
615 doi: 10.1109/CVPR52729.2023.01761. URL [https://doi.org/10.1109/CVPR52729.](https://doi.org/10.1109/CVPR52729.2023.01761)  
616 [2023.01761](https://doi.org/10.1109/CVPR52729.2023.01761).  
617
- 618 Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan  
619 Zhang, Yueze Wang, Zhen Li, Qiyang Yu, Yingli Zhao, Yulong Ao, Xuebin Min, Tao Li,  
620 Boya Wu, Bo Zhao, Bowen Zhang, Lian zi Wang, Guang Liu, Zheqi He, Xi Yang, Jingjing  
621 Liu, Yonghua Lin, Tiejun Huang, and Zhongyuan Wang. Emu3: Next-token prediction is all  
622 you need. *ArXiv*, [abs/2409.18869](https://arxiv.org/abs/2409.18869), 2024. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:272968818)  
623 [CorpusID:272968818](https://api.semanticscholar.org/CorpusID:272968818).  
624
- 625 Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng ming Yin, Shuai  
626 Bai, Xiao Xu, Yilei Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi Wang, An Yang,  
627 Bowen Yu, Chen Cheng, Dayiheng Liu, Deqing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan  
628 Ni, Kai Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng Wang, Shuting Yu, Tingkun  
629 Wen, Wensen Feng, Xiaoxiao Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu, Yuxuan  
630 Cai, and Zenan Liu. Qwen-image technical report, 2025a. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2508.02324)  
631 [2508.02324](https://arxiv.org/abs/2508.02324).  
632
- 633 Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen  
634 Liu, Zhenda Xie, Xingkai Yu, Chong Ruan, and Ping Luo. Janus: Decoupling visual  
635 encoding for unified multimodal understanding and generation. In *IEEE/CVF Confer-*  
636 *ence on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA,*  
637 *June 11-15, 2025*, pp. 12966–12977. Computer Vision Foundation / IEEE, 2025b. URL  
638 [https://openaccess.thecvf.com/content/CVPR2025/html/Wu\\_Janus\\_](https://openaccess.thecvf.com/content/CVPR2025/html/Wu_Janus_Decoupling_Visual_Encoding_for_Unified_Multimodal_Understanding_and_Generation_CVPR_2025_paper.html)  
639 [Decoupling\\_Visual\\_Encoding\\_for\\_Unified\\_Multimodal\\_Understanding\\_](https://openaccess.thecvf.com/content/CVPR2025/html/Wu_Janus_Decoupling_Visual_Encoding_for_Unified_Multimodal_Understanding_and_Generation_CVPR_2025_paper.html)  
640 [and\\_Generation\\_CVPR\\_2025\\_paper.html](https://openaccess.thecvf.com/content/CVPR2025/html/Wu_Janus_Decoupling_Visual_Encoding_for_Unified_Multimodal_Understanding_and_Generation_CVPR_2025_paper.html).  
641
- 642 Chenyuan Wu, Pengfei Zheng, Ruiran Yan, Shitao Xiao, Xin Luo, Yueze Wang, Wanli Li, Xiyan  
643 Jiang, Yexin Liu, Junjie Zhou, Ze Liu, Ziyi Xia, Chaofan Li, Haoge Deng, Jiahao Wang, Kun  
644 Luo, Bo Zhang, Defu Lian, Xinlong Wang, Zhongyuan Wang, Tiejun Huang, and Zheng Liu.  
645 Omnigen2: Exploration to advanced multimodal generation. *arXiv preprint arXiv:2506.18871*,  
646 2025c.  
647
- 648 Yongliang Wu, Zonghui Li, Xinting Hu, Xinyu Ye, Xianfang Zeng, Gang Yu, Wenbo Zhu, Bernt  
649 Schiele, Mingzhuo Yang, and Xu Yang. Kris-bench: Benchmarking next-level intelligent image  
650 editing models. *ArXiv*, [abs/2505.16707](https://arxiv.org/abs/2505.16707), 2025d. URL [https://api.semanticscholar.](https://api.semanticscholar.org/CorpusID:278789151)  
651 [org/CorpusID:278789151](https://api.semanticscholar.org/CorpusID:278789151).

- 648 Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Chao-  
649 fan Li, Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image genera-  
650 tion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025,*  
651 *Nashville, TN, USA, June 11-15, 2025*, pp. 13294–13304. Computer Vision Foundation /  
652 IEEE, 2025. URL [https://openaccess.thecvf.com/content/CVPR2025/html/  
653 Xiao\\_OmniGen\\_Unified\\_Image\\_Generation\\_CVPR\\_2025\\_paper.html](https://openaccess.thecvf.com/content/CVPR2025/html/Xiao_OmniGen_Unified_Image_Generation_CVPR_2025_paper.html).
- 654 Siwei Yang, Mude Hui, Bingchen Zhao, Yuyin Zhou, Nataniel Ruiz, and Cihang Xie. Complex-  
655 edit: Cot-like instruction generation for complexity-controllable image editing benchmark. *CoRR*,  
656 abs/2504.13143, 2025. doi: 10.48550/ARXIV.2504.13143. URL [https://doi.org/10.  
657 48550/arXiv.2504.13143](https://doi.org/10.48550/arXiv.2504.13143).
- 658 Qifan Yu, Wei Chow, Zhongqi Yue, Kaihang Pan, Yang Wu, Xiaoyang Wan, Juncheng Li, Siliang  
659 Tang, Hanwang Zhang, and Yueting Zhuang. Anyedit: Mastering unified high-quality image edit-  
660 ing for any idea. *ArXiv*, abs/2411.15738, 2024. URL [https://api.semanticscholar.  
661 org/CorpusID:274233770](https://api.semanticscholar.org/CorpusID:274233770).
- 662 Daoan Zhang, Che Jiang, Ruoshi Xu, Biaoxiang Chen, Zijian Jin, Yutian Lu, Jianguo Zhang, Liang  
663 Yong, Jiebo Luo, and Shengda Luo. Worldgenbench: A world-knowledge-integrated benchmark  
664 for reasoning-driven text-to-image generation. *ArXiv*, abs/2505.01490, 2025. URL [https:  
665 //api.semanticscholar.org/CorpusID:278327313](https://api.semanticscholar.org/CorpusID:278327313).
- 666 Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. Magicbrush: A manually  
667 annotated dataset for instruction-guided image editing. In Alice Oh, Tristan Naumann,  
668 Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in*  
669 *Neural Information Processing Systems 36: Annual Conference on Neural Information*  
670 *Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*  
671 *2023*, 2023a. URL [http://papers.nips.cc/paper\\_files/paper/2023/  
672 hash/64008fa30cba9b4d1ab1bd3bd3d57d61-Abstract-Datasets\\_and\\_  
673 Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/64008fa30cba9b4d1ab1bd3bd3d57d61-Abstract-Datasets_and_Benchmarks.html).
- 674 Xinlu Zhang, Yujie Lu, Weizhi Wang, An Yan, Jun Yan, Lianke Qin, Heng Wang, Xifeng Yan,  
675 William Yang Wang, and Linda Ruth Petzold. Gpt-4v(ision) as a generalist evaluator for vision-  
676 language tasks. *CoRR*, abs/2311.01361, 2023b. doi: 10.48550/ARXIV.2311.01361. URL  
677 <https://doi.org/10.48550/arXiv.2311.01361>.
- 678 Xiangyu Zhao, Peiyuan Zhang, Kexian Tang, Hao Li, Zicheng Zhang, Guangtao Zhai, Junchi  
679 Yan, Hua Yang, Xue Yang, and Haodong Duan. Envisioning beyond the pixels: Benchmark-  
680 ing reasoning-informed visual editing. *ArXiv*, abs/2504.02826, 2025. URL [https://api.  
681 semanticscholar.org/CorpusID:277510499](https://api.semanticscholar.org/CorpusID:277510499).
- 682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A APPENDIX

### A.1 OVERVIEW OF THE APPENDIX

This appendix supplements the proposed InEdit-Bench with details excluded from the main paper due to space constraints.

The appendix is organized as follows:

- Sec. A.2: More Detailed Evaluation Results.
- Sec. A.3: Data Source of InEdit-Bench.
- Sec. A.4: Limitations.
- Sec. A.5: Utilization of Large Language Models.
- Sec. A.6: Representative Examples of InEdit-Bench Subtasks.
- Sec. A.7: Detailed Outputs of Evaluated Models.
- Sec. A.8: Design of the Prompt.

### A.2 MORE DETAILED EVALUATION RESULTS

In this section, we present a more detailed evaluation of the models, offering further analysis of their capabilities. This includes:

- (1) The specific scores of 14 models across four fundamental task categories and 16 sub-tasks.
- (2) The accuracy of these 14 models across the 16 sub-tasks.

#### A.2.1 SCORES OF MODELS ACROSS 4 TASKS AND 16 SUB-TASKS

Tab. 2 and Fig. 7 show the specific scores of 14 models across four fundamental tasks and 16 sub-tasks, respectively. Compared to open-source models, proprietary models exhibit significantly more balanced performance. In all tasks, the models particularly excel in the perceptual quality assessment dimension, demonstrating their ability to generate natural, smooth, and high-quality images, avoiding common issues such as distortion and blurring. This result indicates that current models effectively address challenges related to image quality when generating visual content.

Most models score slightly lower in the appearance consistency dimension compared to the perceptual quality dimension, yet still demonstrate considerable capability. However, some models have still failed to effectively adapt to the new paradigm of intermediate logic path editing, resulting in poor performance in the appearance consistency dimension. For example, models like Emu1 and Omnigen encounter significant obstacles in this dimension, with performance far below that of other models.

There are significant variations in performance across models in terms of semantic consistency and logical consistency. Except for GPT-Image-1 and Nano-Banana, other models show significant imbalances in these two dimensions, with a noticeable drop in scores. Notably, models like Emu1 and Omnigen score almost zero in both semantic consistency and logical consistency, highlighting the limitations of current models in handling complex logical relationships.

Among the open-source models, all except Qwen-Image-Edit, Bagel, and Bagel-Think show relatively poor performance. Among them, a few models, such as Omnigen2 and Step1X-Edit(v1.1), show slight improvements in certain sub-tasks, achieving some scores. However, these advancements are not applicable to a broader range of tasks. Overall, the performance of open-source models still lags significantly behind that of proprietary models, with notable gaps in multiple key dimensions.

#### A.2.2 ACCURACY OF MODELS ACROSS 16 SUB-TASKS

Tab. 3 shows the accuracy scores of each model across four basic categories and 16 sub-tasks. The effective scores are primarily concentrated in GPT-Image-1 and Nano-Banana. Both models perform well in multiple sub-tasks. Although GPT-Image-1 outperforms Nano-Banana overall, Nano-

Table 2: The specific scores of the models across four fundamental tasks, with metrics including Appearance Consistency (AC), Perceptual Quality (PQ), Semantic Consistency (SC), Logical Coherence (LC), Scientific Plausibility (SP). The performance of open-source and proprietary models is separately marked with the best performance in **bold**, and the second best underlined.

Metric	Proprietary Models				Open-Source Models										
	GPT-Image-1	Nano-Banana	Flux-Komext-pro	Doubao-SeedEdit-3.0	Owen-Image-Edit	Emu1	Emu2	Bagel	Bagel-Think	OmniGen	OmniGen2	Step1X-Edit (v1.0)	Step1X-Edit (v1.1)	InstructPix2Pix	
<i>State Transition</i>	AC	<b>95.4</b>	<u>79.6</u>	56.1	42.7	<b>53.1</b>	3.1	30.1	38.3	<u>48.5</u>	5.6	31.6	12.8	21.9	24.5
	PQ	92.3	<u>94.4</u>	<b>94.8</b>	69.4	<u>81.6</u>	44.9	<b>83.3</b>	69.4	78.1	29.1	80.1	42.7	47.4	73.0
	SC	<b>72.4</b>	<u>58.7</u>	31.6	21.9	<b>24.5</b>	1.5	8.7	<u>24.0</u>	19.9	2.0	19.4	7.7	10.7	4.6
	LC	<b>64.3</b>	<u>48.0</u>	24.5	19.4	19.9	1.0	11.7	<u>21.4</u>	20.9	5.1	<b>22.4</b>	5.6	13.8	12.2
	Avg	<b>81.1</b>	<u>70.2</u>	51.8	38.4	<b>44.8</b>	12.6	33.5	38.3	41.8	10.5	38.4	17.2	23.5	28.6
<i>Temporal Sequence</i>	AC	<b>89.4</b>	<u>85.2</u>	64.4	49.2	<b>62.1</b>	6.8	35.0	47.7	<u>55.7</u>	9.5	47.7	16.3	35.2	45.1
	PQ	<b>89.8</b>	<u>87.1</u>	83.7	69.3	<u>83.3</u>	45.5	<b>87.3</b>	64.4	75.0	40.2	79.5	37.5	50.8	78.5
	SC	<b>71.6</b>	<u>64.4</u>	33.7	23.1	<u>28.4</u>	2.7	1.5	<b>29.2</b>	20.8	6.1	21.2	6.8	16.7	5.4
	LC	<b>74.2</b>	<u>60.2</u>	26.9	27.3	<b>34.5</b>	4.9	14.4	<u>29.2</u>	27.7	9.5	25.4	8.3	26.1	19.3
	Avg	<b>81.3</b>	<u>74.2</u>	52.2	42.2	<b>52.1</b>	15.0	34.6	42.6	<u>44.8</u>	16.3	43.5	17.2	32.2	37.1
<i>Dynamic Process</i>	AC	<u>91.9</u>	<b>92.7</b>	70.8	44.6	<b>71.9</b>	5.0	34.6	51.5	<u>58.1</u>	12.3	46.9	20.4	41.9	30.4
	PQ	<u>94.6</u>	<b>97.3</b>	93.5	70.0	<u>85.8</u>	55.0	<b>88.5</b>	64.6	<u>76.2</u>	36.9	79.2	48.0	65.4	71.5
	SC	<b>71.9</b>	<u>62.7</u>	36.2	23.1	<u>28.8</u>	3.8	8.1	<b>31.2</b>	<u>29.6</u>	7.7	23.1	9.2	20.4	3.8
	LC	<b>70.4</b>	<u>64.6</u>	36.9	21.5	<u>30.4</u>	4.2	18.1	30.0	<b>35.4</b>	9.2	28.5	10.8	28.8	11.5
	SP	<u>70.4</u>	<b>71.9</b>	43.8	38.5	<b>47.7</b>	4.6	21.2	34.2	<u>35.4</u>	14.6	31.5	16.5	31.2	14.6
Avg	<b>79.8</b>	<u>77.8</u>	56.2	39.5	<b>52.9</b>	14.5	34.1	42.3	<u>46.9</u>	16.2	41.8	21.0	37.5	26.4	
<i>Scientific Simulation</i>	AC	<b>94.6</b>	<u>87.0</u>	66.3	33.7	<b>55.4</b>	5.4	30.4	43.5	<u>48.9</u>	7.6	37.0	8.7	39.1	29.3
	PQ	<b>93.5</b>	<u>90.2</u>	80.4	70.7	73.9	47.8	<u>75.0</u>	66.3	<b>80.4</b>	33.7	73.9	42.4	50.0	<u>75.0</u>
	SC	<b>72.8</b>	<u>68.5</u>	33.7	20.7	<b>26.1</b>	0.0	8.7	<u>25.0</u>	<b>26.1</b>	4.3	20.7	3.3	21.7	3.3
	LC	<b>78.3</b>	<u>73.9</u>	32.6	17.4	<u>28.3</u>	2.2	16.3	27.2	<b>31.5</b>	4.3	17.4	4.3	25.0	4.3
	SP	<u>73.9</u>	<b>78.3</b>	43.5	25.0	<b>37.0</b>	6.5	26.1	<u>33.7</u>	30.4	12.0	22.8	10.9	33.7	9.8
Avg	<b>82.6</b>	<u>79.6</u>	51.3	33.5	<b>44.1</b>	12.4	31.3	39.1	<u>43.5</u>	12.4	34.3	13.9	33.9	24.3	

Banana still has an advantage in certain sub-tasks. In the case of open-source models, all models have an accuracy of 0% in the state transition and scientific simulation category tasks, and in tasks from other categories, only a few models show slight improvements in their scores. Overall, even the most advanced models achieve only 16.75% accuracy, with more than half of the models scoring 0%. This indicates that current models are still in the early stages of solving intermediate logic path editing tasks, far from meeting the requirements for practical application.

### A.3 DATA SOURCE OF INEDIT-BENCH

Input images for the InEdit-Bench dataset are primarily sourced from the following categories:

- (1) Images generated by image generation models.
- (2) Images derived from existing datasets and benchmarks.
- (3) Images collected from the internet under permissive licenses.

### A.4 LIMITATIONS

This study aims to establish a pioneering benchmark for intermediate logical reasoning and multi-step editing tasks. However, as an initial exploration, the current benchmark still has several aspects that require improvement. We openly acknowledge its potential limitations, such as the dataset’s insufficient scale to cover all complex scenarios and the task categorization that may not exhaust all possibilities. Future work will focus on addressing these issues to build a more comprehensive and robust benchmark.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

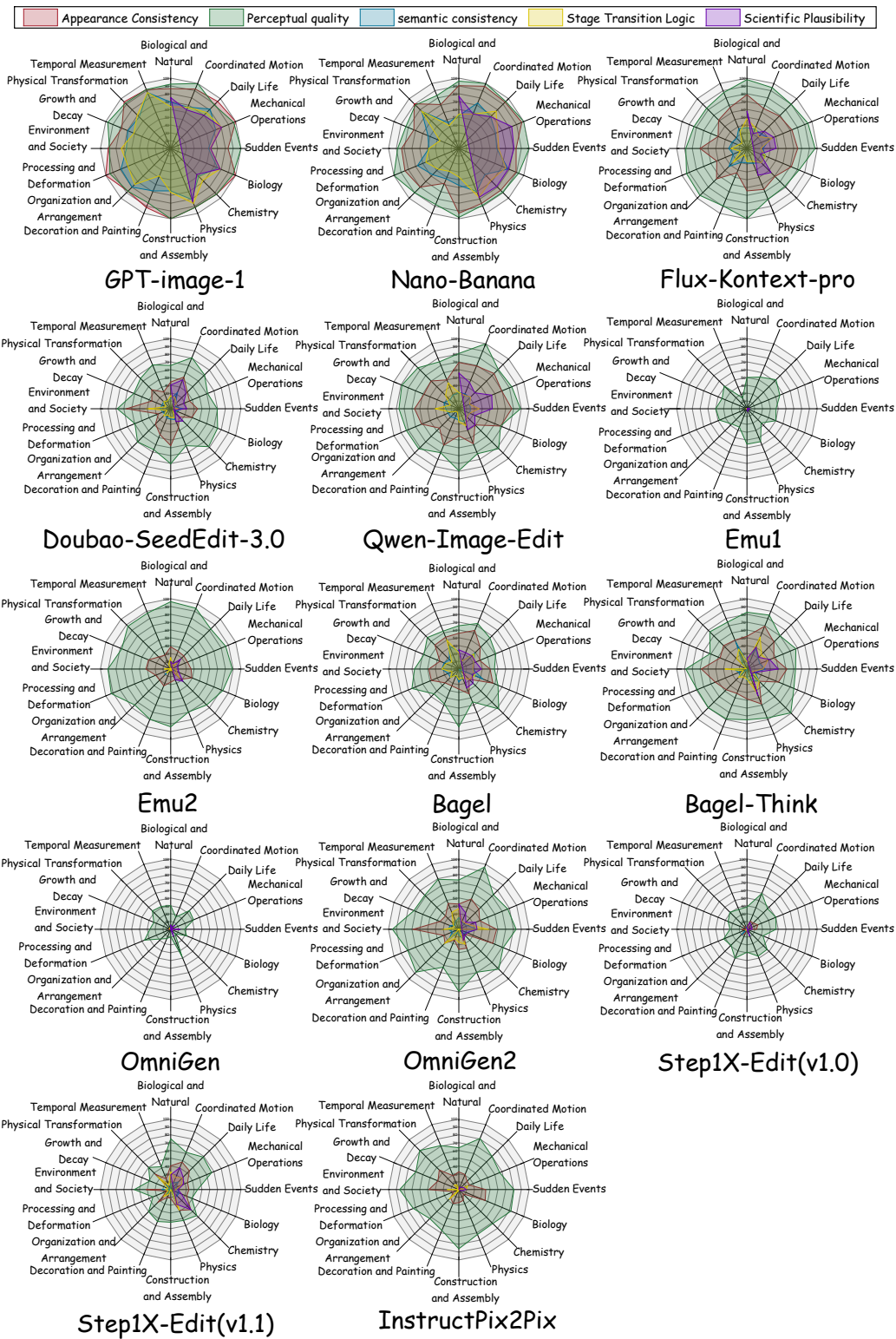


Figure 7: The average scores for 14 models across 16 sub-tasks.

Table 3: Accuracy performance of different models across 16 sub-tasks, including State Transition: Construction and Assembly (CA), Decoration and Painting (DP), Organization and Arrangement (OA), Processing and Deformation (PD). Temporal Sequence: Environment and Society (ES), Growth and Decay (GD), Physical Transformation (PT), Temporal Measurement (TM). Dynamic Process: Biology and Nature (BN), Coordinated Motion (CM), Daily Life (DL), Mechanical Operations (MO), Sudden Events (SE). Scientific Simulation: Biology (BI), Chemistry (CH), Physics (PH). The performance of open-source and proprietary models is separately marked, with the best performance in **bold** and the second-best performance underlined.

SubTasks	Proprietary Models				Open-Source Models										
	GPT-Image-1	Nano-Banana	Flux-Kontext-pro	Doubao-SeedEdit-3.0	Qwen-Image-Edit	Emu1	Emu2	Bagel	Bagel-Think	OmniGen	OmniGen2	StepIX-Edit (v1.0)	StepIX-Edit (v1.1)	InstructPix2Pix	
<i>State Transition</i>	CA	<b>14.29</b>	<u>7.14</u>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	DP	<b>8.33</b>	<b>8.33</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	OA	<b>40.00</b>	<u>10.00</u>	<u>10.00</u>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	PD	<u>7.69</u>	<b>15.38</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Avg	<b>16.33</b>	<u>10.20</u>	2.04	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
<i>Temporal Sequence</i>	ES	<u>10.53</u>	<b>15.79</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>5.26</b>	0.00	0.00	0.00
	GD	<b>20.00</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	PT	<u>12.00</u>	<b>40.00</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	TM	<b>57.14</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>14.29</b>	0.00	0.00	0.00	0.00	0.00
	Avg	<u>18.18</u>	<b>19.70</b>	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.52</b>	0.00	<b>1.52</b>	0.00	0.00	0.00
<i>Dynamic Process</i>	BN	<b>15.38</b>	<u>7.69</u>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	CM	0.00	<b>7.69</b>	0.00	0.00	<b>7.69</b>	0.00	0.00	0.00	<b>7.69</b>	0.00	0.00	0.00	0.00	0.00
	DL	<b>28.57</b>	<u>9.52</u>	4.76	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	MO	<b>22.22</b>	<b>22.22</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	SE	<b>11.11</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Avg	<b>16.92</b>	<u>9.23</u>	1.54	0.00	<b>1.54</b>	0.00	0.00	0.00	<b>1.54</b>	0.00	0.00	0.00	0.00	0.00	
<i>Scientific Simulation</i>	BI	0.00	<b>28.57</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	CH	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	PH	<b>33.33</b>	<u>11.11</u>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Avg	<b>13.04</b>	<b>13.04</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
<b>Overall Accuracy</b>	<b>16.75</b>	<u>13.30</u>	0.99	0.00	<u>0.49</u>	0.00	0.00	0.00	<b>0.99</b>	0.00	<u>0.49</u>	0.00	0.00	0.00	

## A.5 UTILIZATION OF LARGE LANGUAGE MODELS

The core methodology of this study was independently designed, while the assistance of LLMs enhanced the efficiency and completeness of the research across several stages. Specifically: (1) Evaluation dataset construction: during the process of building the evaluation dataset, we employed LLMs to assist in conceptualizing instance scenarios, thereby improving the comprehensiveness of scenario coverage. (2) Knowledge checklist design: the checklist incorporates the key mechanisms and features of intermediate logic paths, and LLMs were leveraged to aid its design and refinement, ensuring both scientific rigor and validity. By integrating LLMs in these stages, we were able to exploit their advanced language understanding capabilities while further optimizing the research workflow, making the overall study more comprehensive and robust.

## A.6 REPRESENTATIVE EXAMPLES OF INEDIT-BENCH SUBTASKS

In this section, we present representative example images from the 16 subtasks in InEdit-Bench, with each subtask corresponding to a distinct testing scenario. Fig. 8–11 illustrate examples from the four task categories: 4 subtasks of state transition, 4 subtasks of temporal sequence, 5 subtasks of dynamic process, and 3 subtasks of scientific simulation.

## A.7 DETAILED OUTPUTS OF EVALUATED MODELS

Some of the evaluated model outputs from our InEdit-Bench benchmark are shown in Fig. 12–24, providing a more intuitive understanding of the performance of the tested models.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

## A.8 DESIGN OF THE PROMPT

In this section, we specifically present the instruction prompts and evaluation prompts used for intermediate logic path editing.

### A.8.1 EDIT PROMPT

Fig. 25 shows the instructions we used to generate intermediate logic path editing results. For each instruction, the overall structure is as follows: first, we briefly introduce the starting and ending state goals and request the generation of the logical transition process in between. Then, we standardize the output format, requiring the output image to be divided into N grids, with each grid representing a node. Finally, to guide the tested model in clearly presenting the intermediate process rather than focusing on redundant node information, we add prompts for key nodes. For State Transition category tasks, we require the model to treat each step of the intermediate process as a key node. For Temporal Sequence category tasks, we require the model to divide the entire intermediate process into equal time intervals. For Dynamic Process and Scientific Simulation category tasks, we use a large multimodal model to help briefly define some key nodes. Additionally, for the Path Understanding section, we manually annotated the sequence that the intermediate logic path should follow.

### A.8.2 EVALUATION PROMPT

Fig. 27–32 specifically show the prompts we used for evaluation. Additionally, in the Scientific Plausibility evaluation dimension, there is a Knowledge Checklist that includes key features or intrinsic mechanisms of the intermediate process. Fig. 26 presents a sample instance, where each sample contains 2-4 inspection items along with corresponding explanations, guiding the model to better understand the evaluation principles through the item descriptions.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

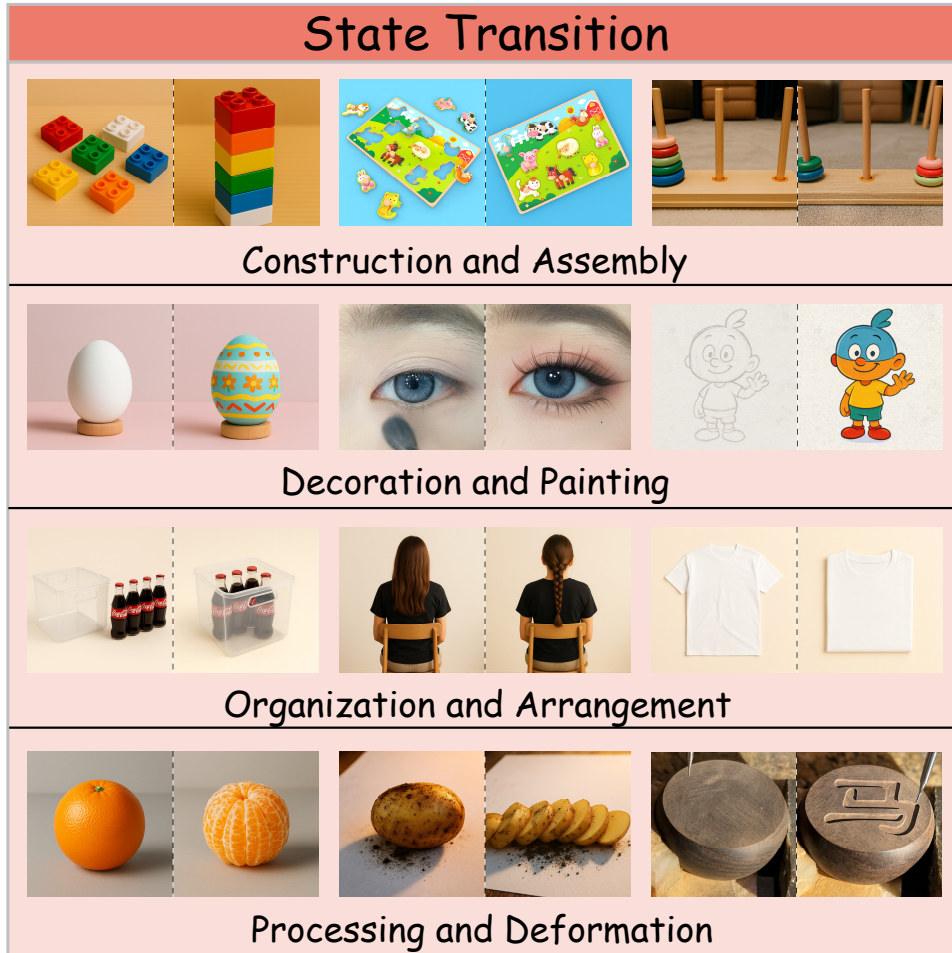


Figure 8: Representative examples of state transition.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

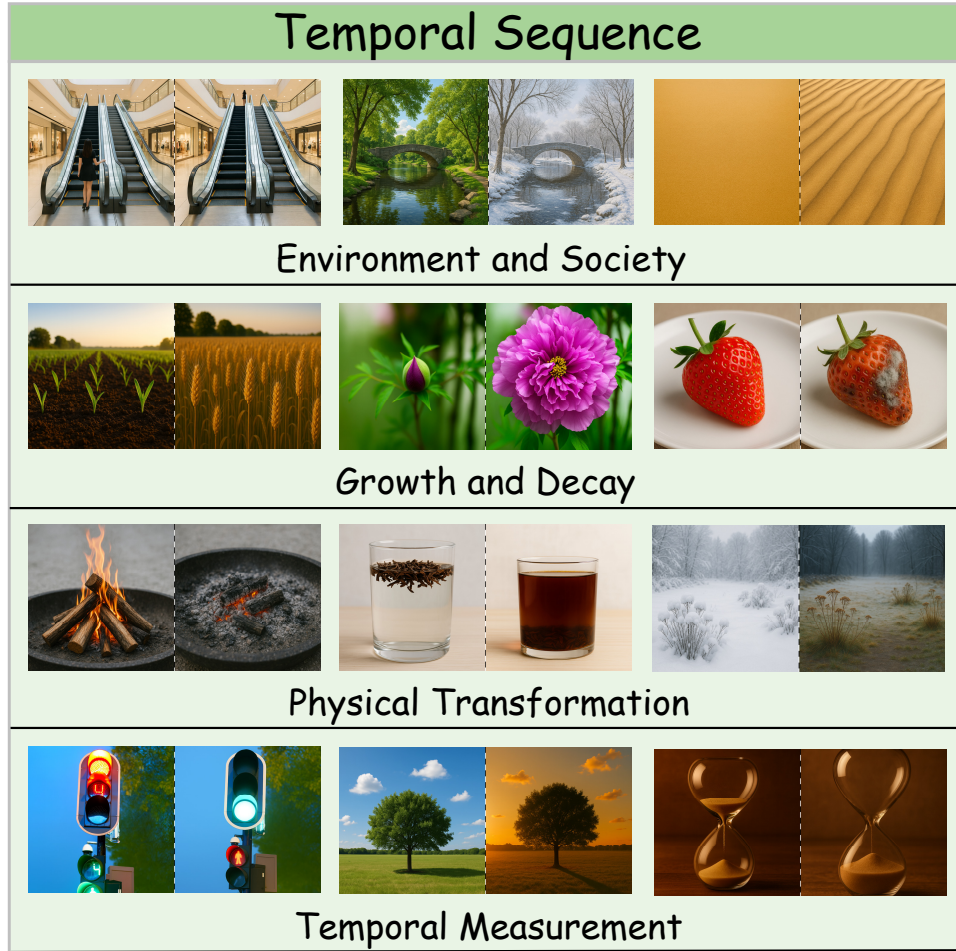


Figure 9: Representative examples of temporal sequence.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

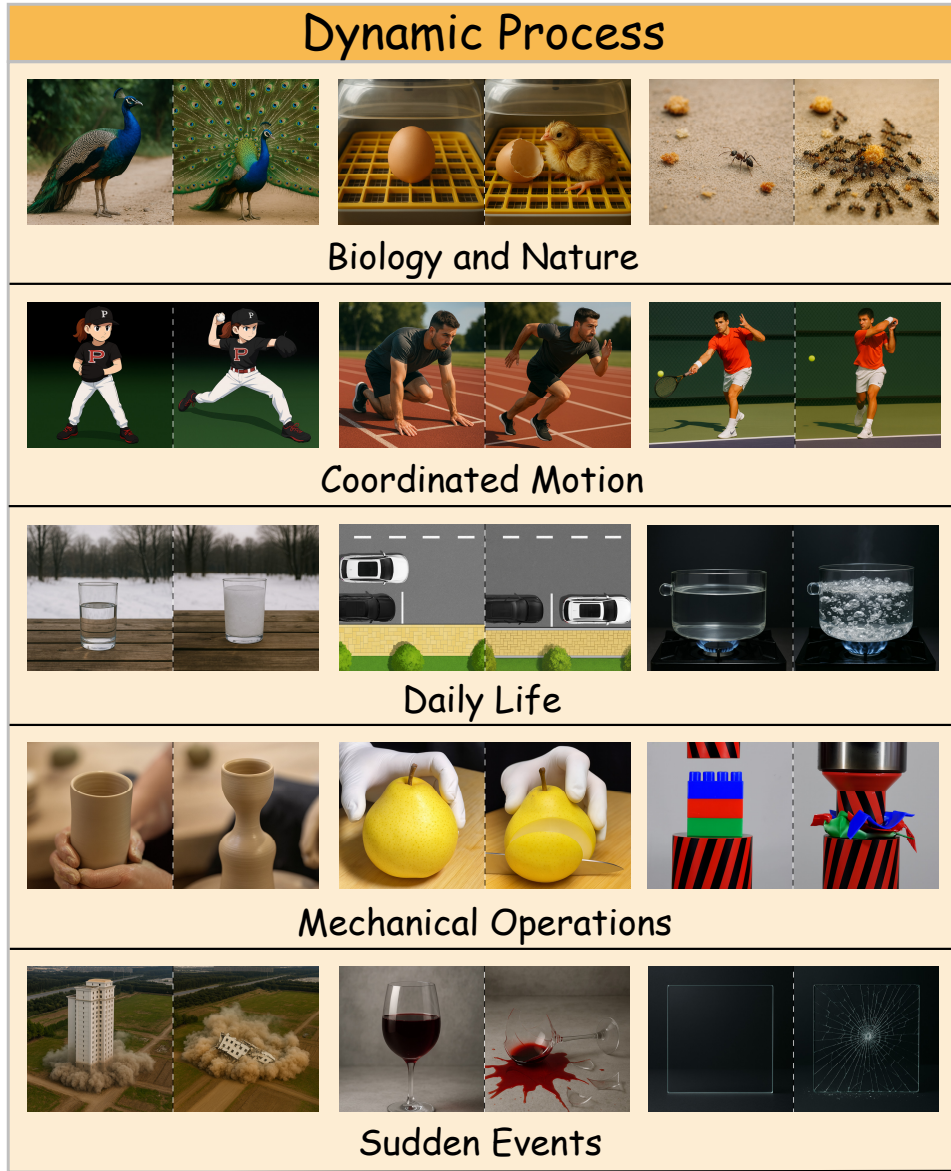


Figure 10: Representative examples of dynamic process.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

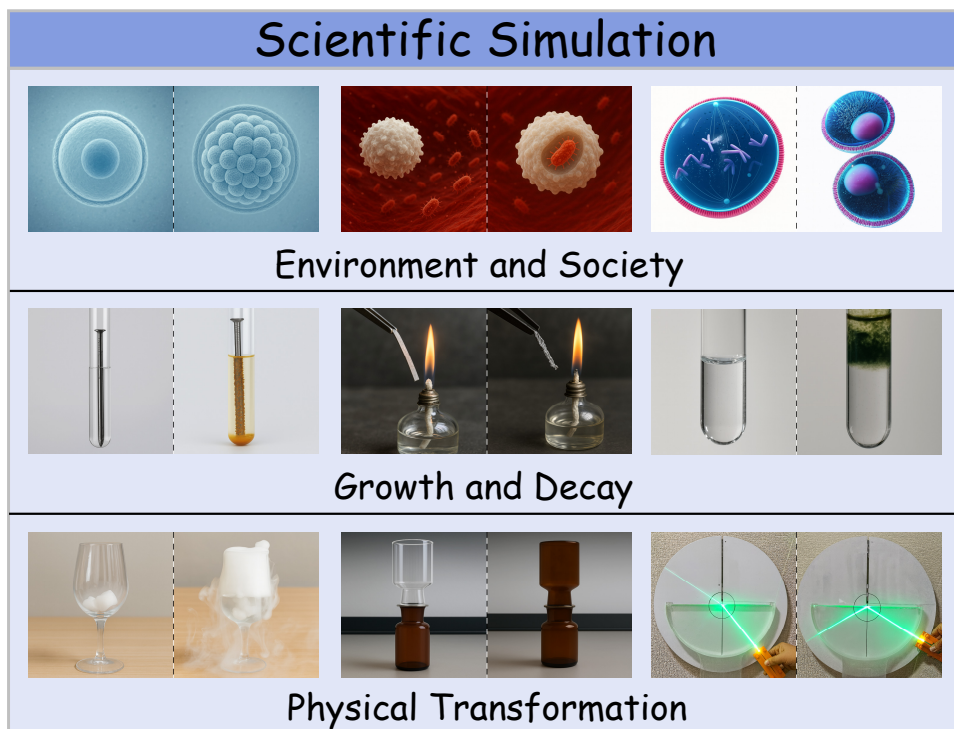


Figure 11: Representative examples of scientific simulation.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241



Figure 12: State Transition Outputs - Part1.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295



Figure 13: State Transition Outputs - Part2.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

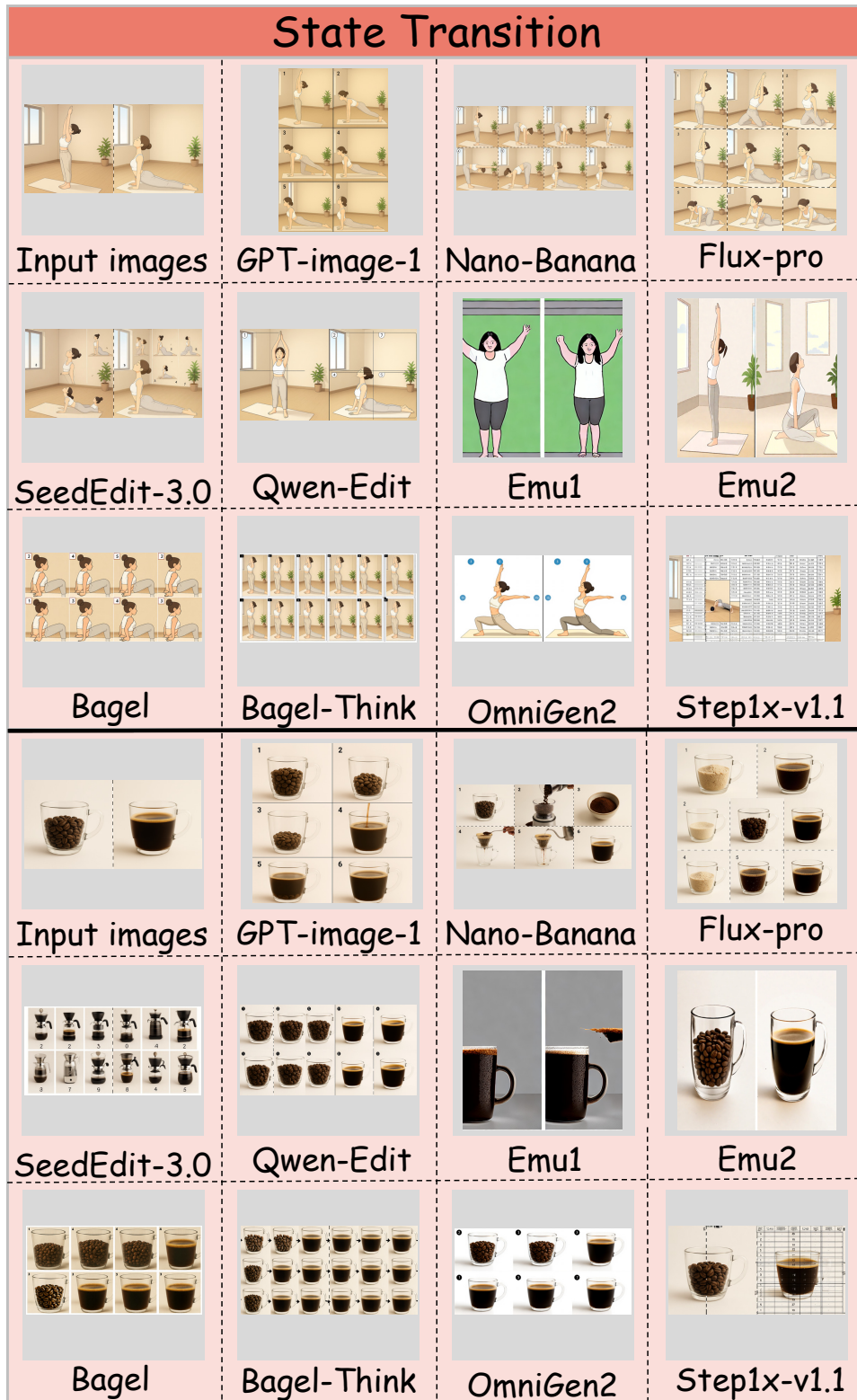


Figure 14: State Transition Outputs - Part3.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

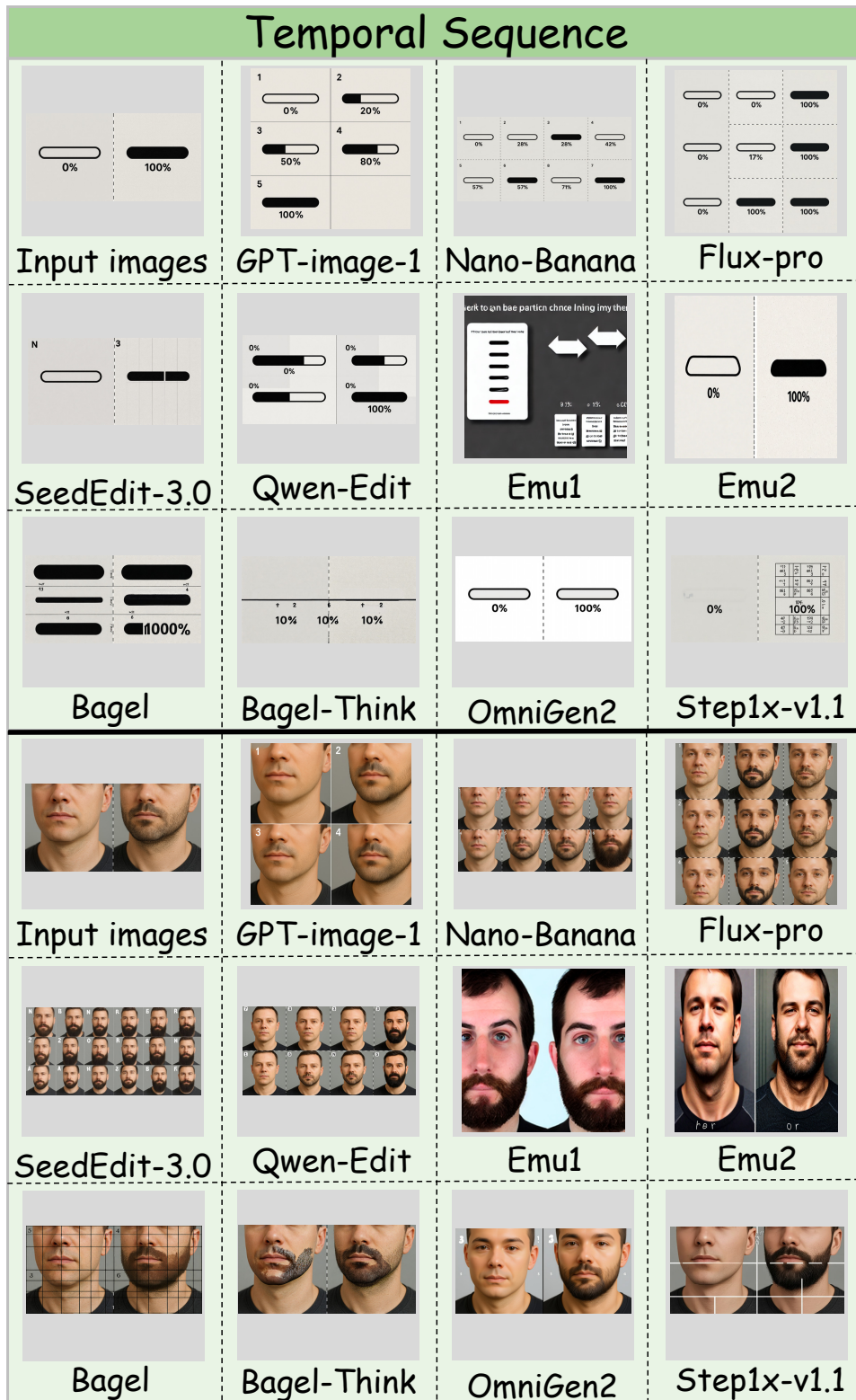


Figure 15: Temporal Sequence Outputs - Part1.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

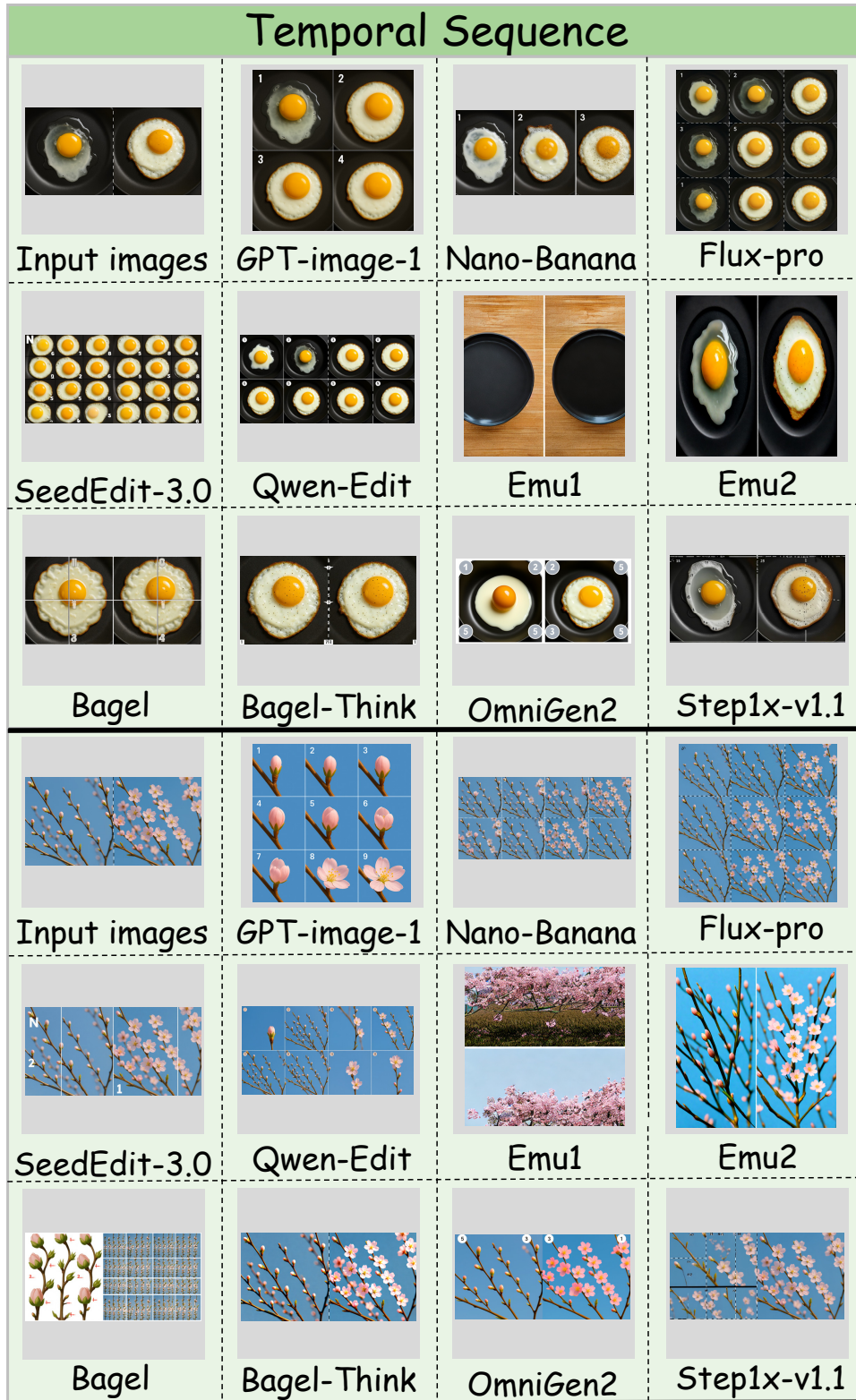


Figure 16: Temporal Sequence Outputs - Part2.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

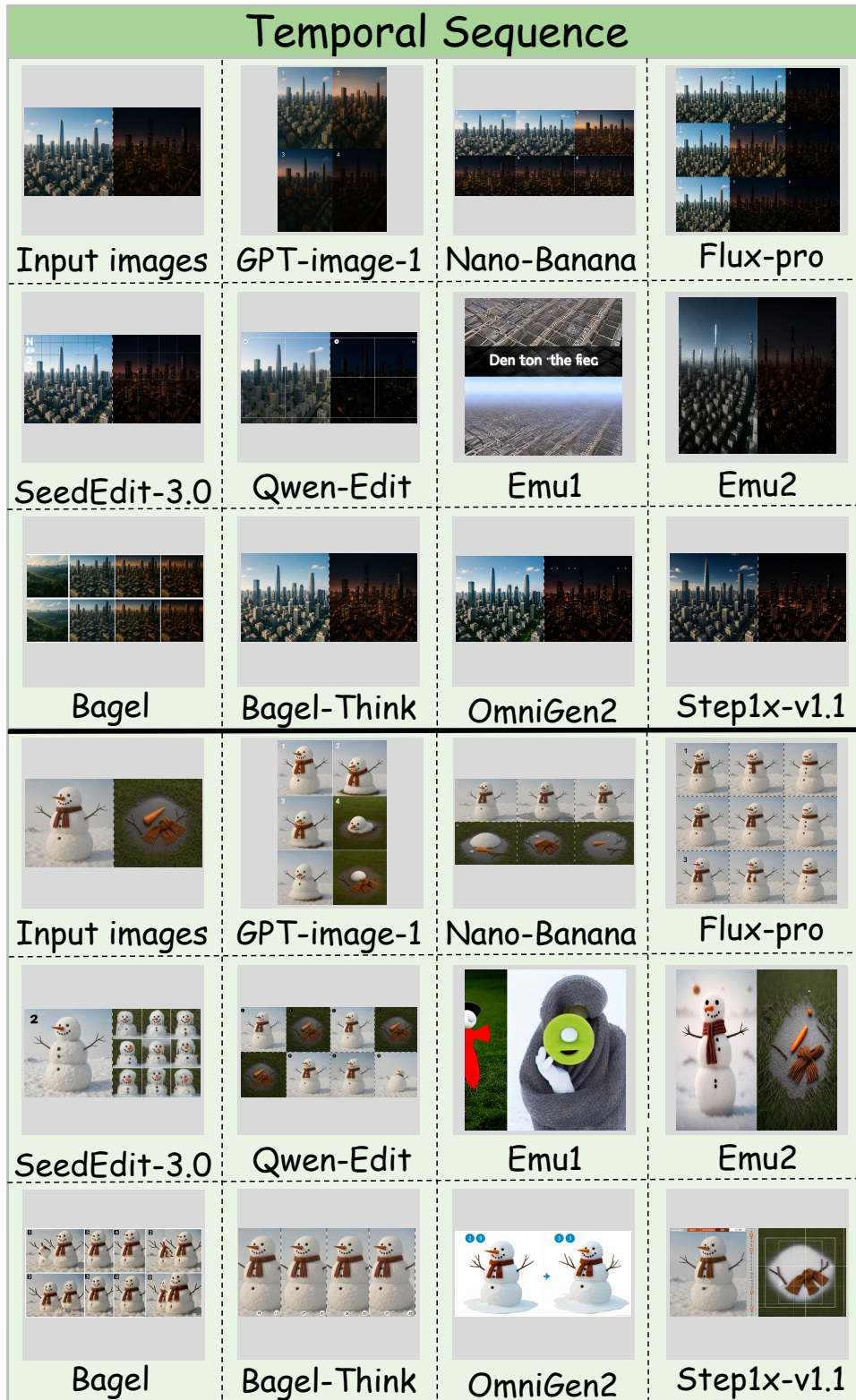


Figure 17: Temporal Sequence Outputs - Part3.

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

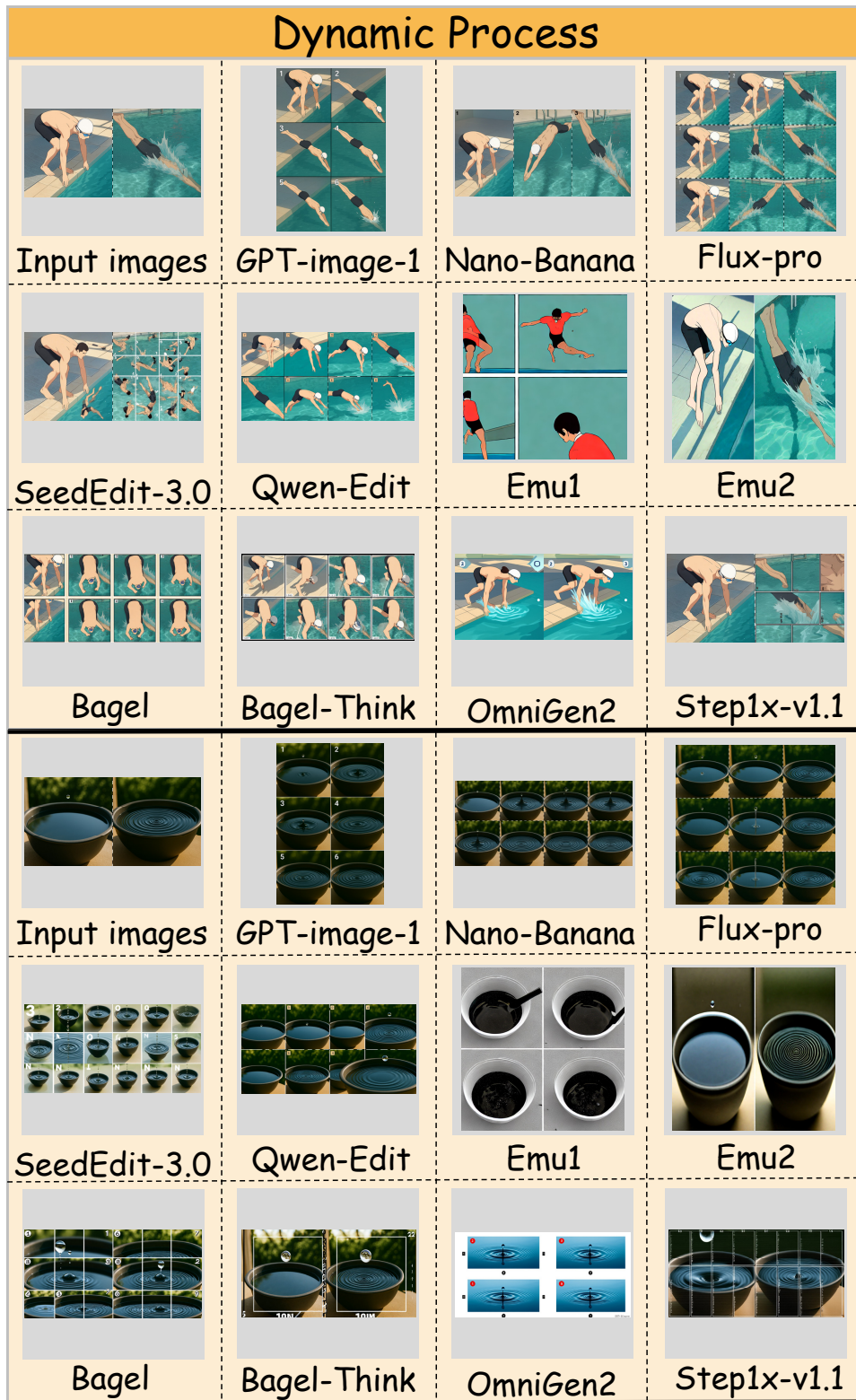


Figure 18: Dynamic Process Outputs - Part1.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

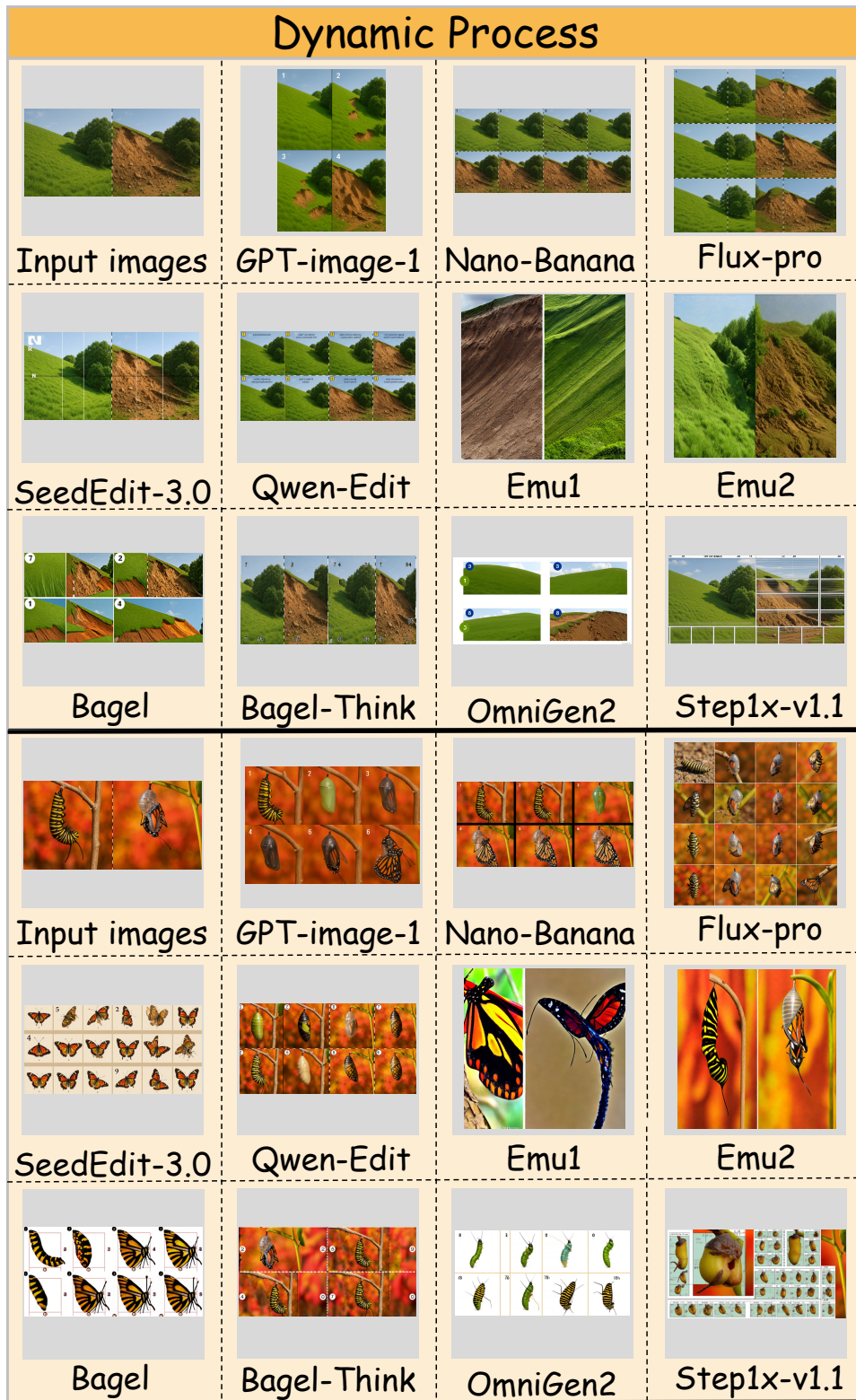


Figure 19: Dynamic Process Outputs - Part2.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

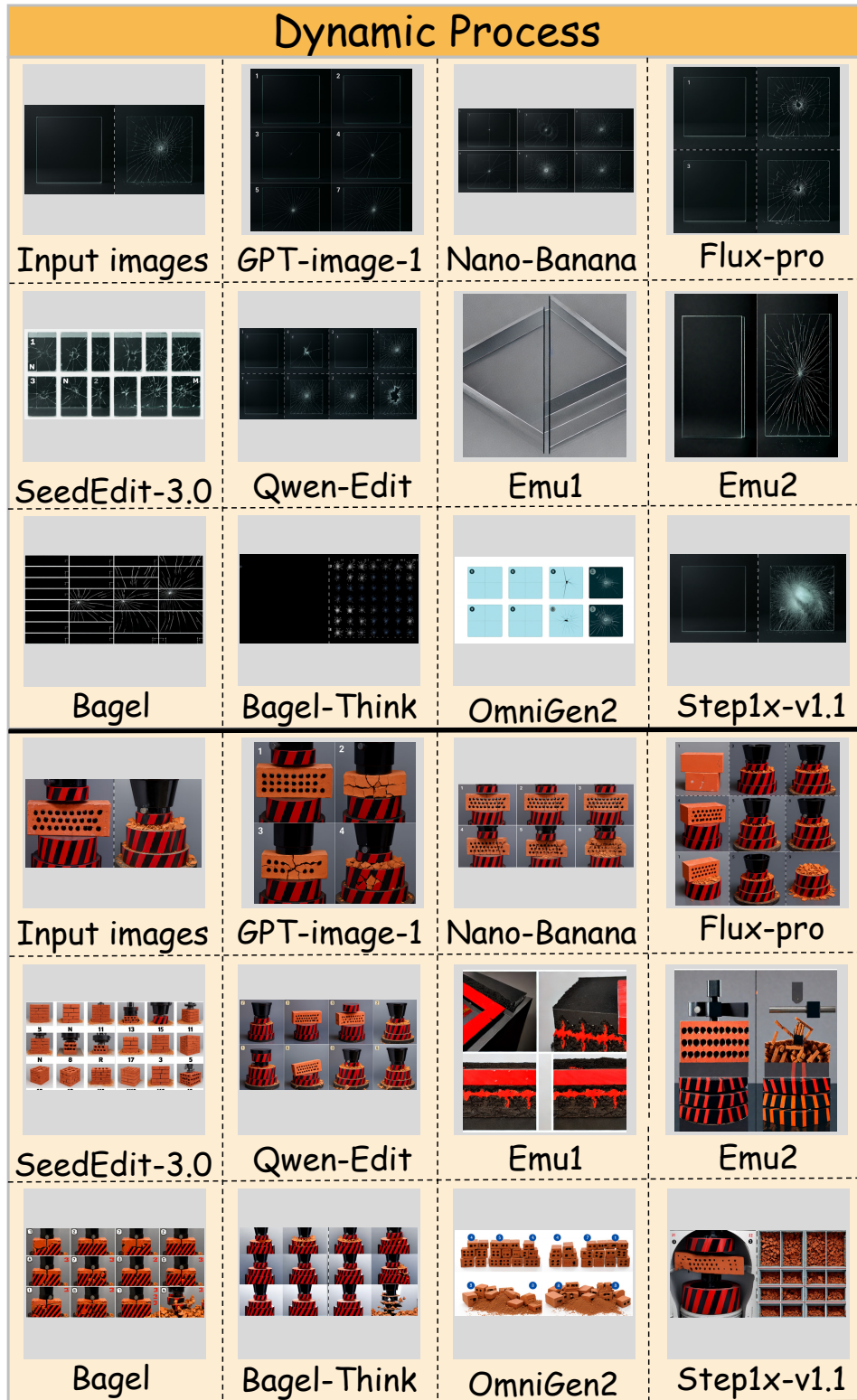


Figure 20: Dynamic Process Outputs - Part3.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

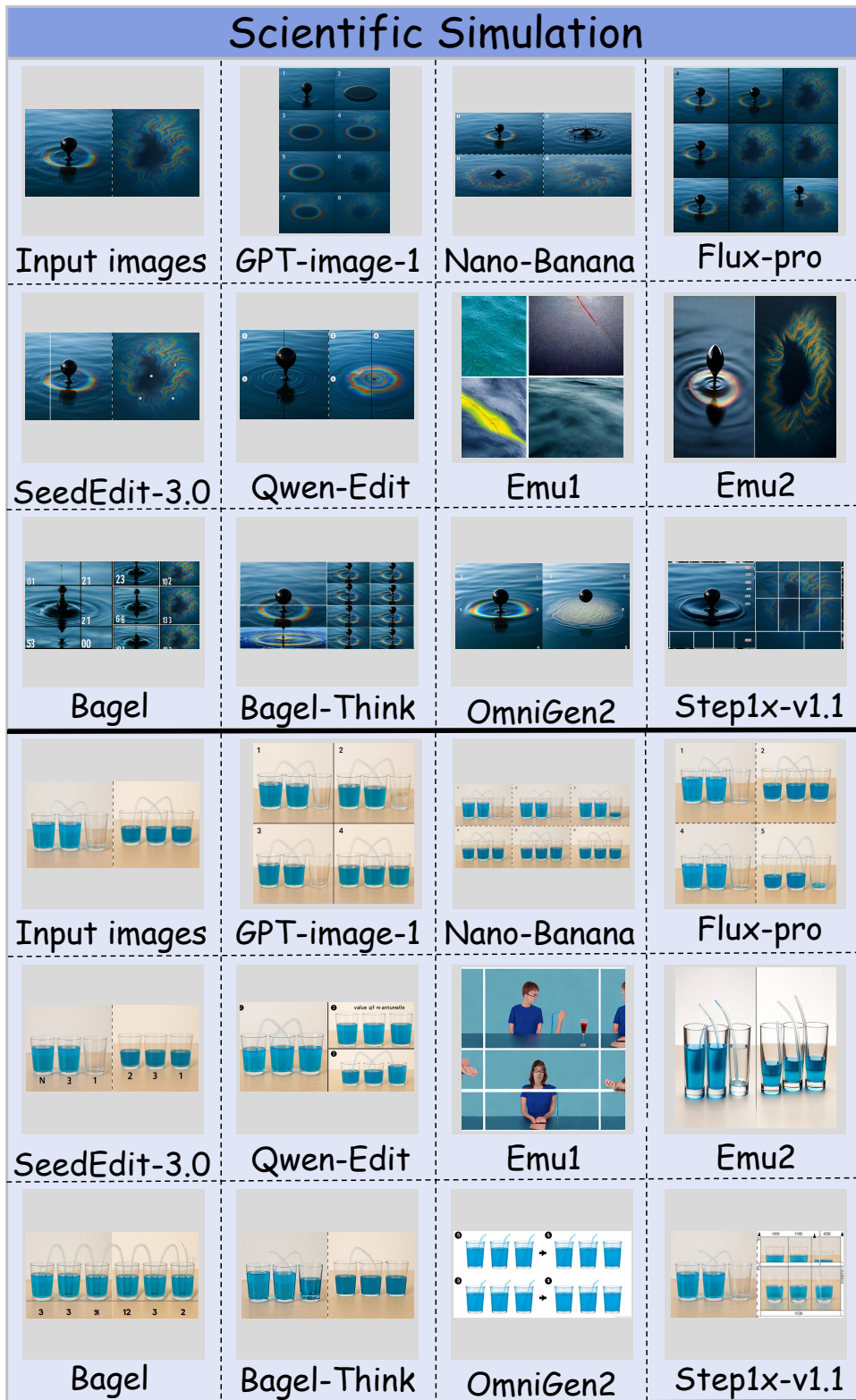


Figure 21: Scientific Simulation Outputs - Part1.

1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781

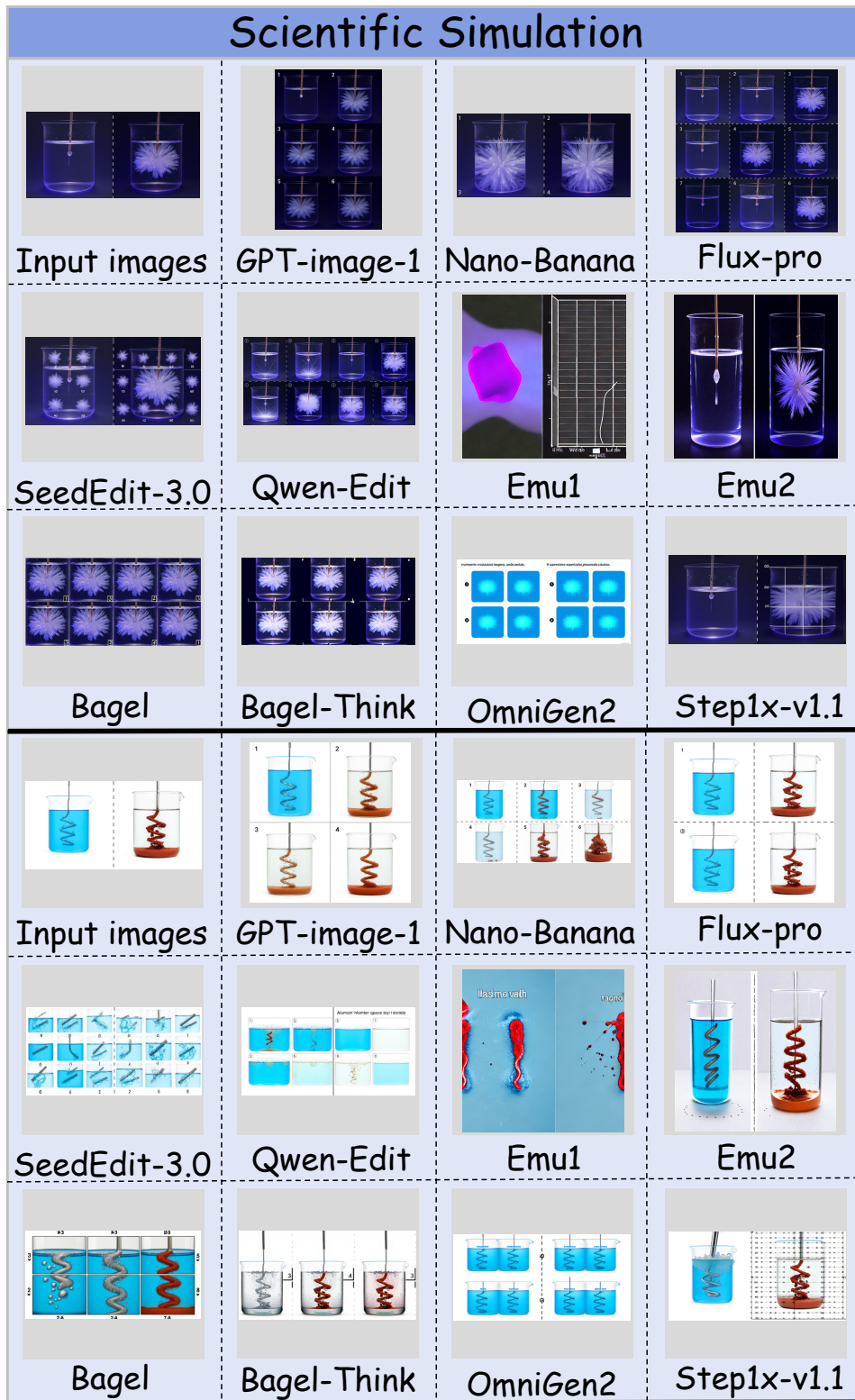


Figure 22: Scientific Simulation Outputs - Part2.

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

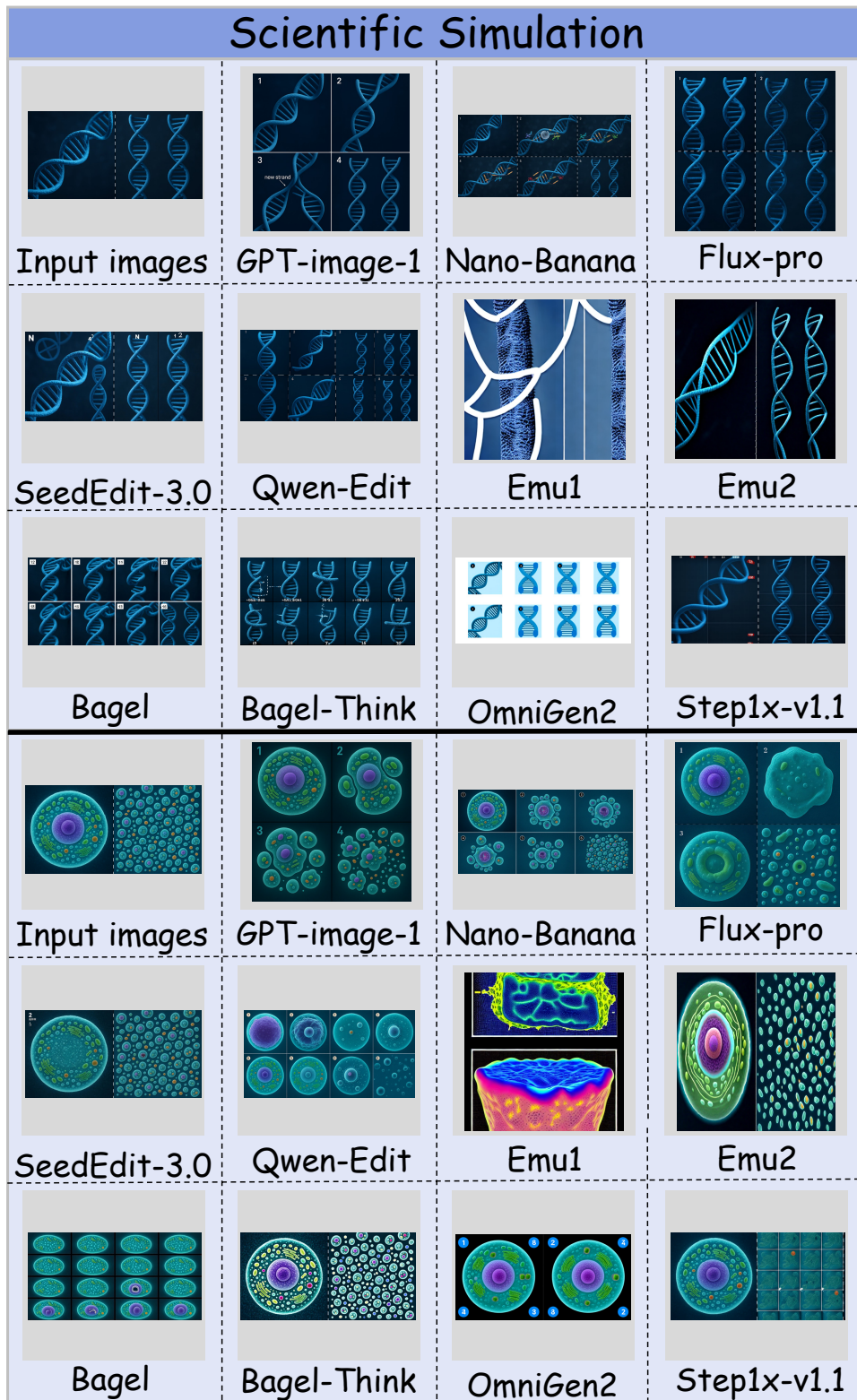


Figure 23: Scientific Simulation Outputs - Part3.

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

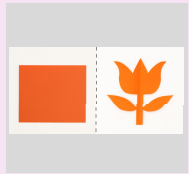
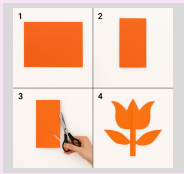
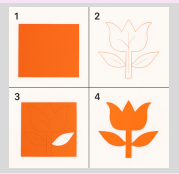
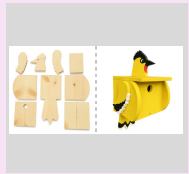
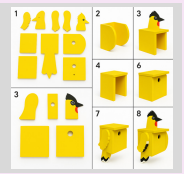
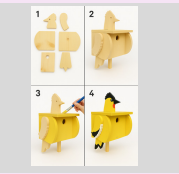
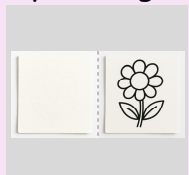

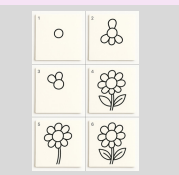



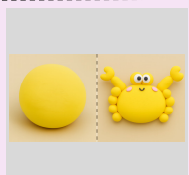
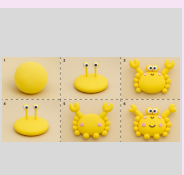
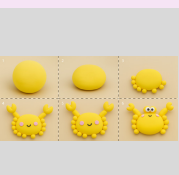
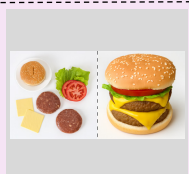

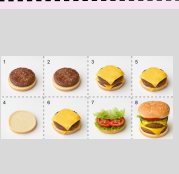
Process Plausibility				
			<b>Path_A:</b> Fold (center) → Cut (half flower) → Unfold. <b>Path_B:</b> Draw outline → Cut along line → Done.	
			<b>Path_A:</b> Paint boards → Assemble. <b>Path_B:</b> Assemble birdhouse → Paint.	
			<b>Path_A:</b> Outline → Petals → Leaves & Stem → Center & Details. <b>Path_B:</b> Center → Petals → Stem → Leaves.	
			<b>Path_A:</b> Insert bouquet stepwise. <b>Path_B:</b> Finally insert all bouquets together.	
			<b>Path_A:</b> torso → eyes → claws → legs → blush → mouth. <b>Path_B:</b> torso → legs → claws → blush → mouth → eyes.	
			<b>Path_A:</b> Prepare middle → Add bottom → Add top. <b>Path_B:</b> Bottom bun → Stack ingredients → Top bun.	

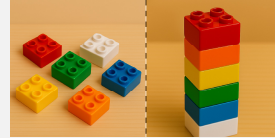
Figure 24: Process Plausibility Outputs.

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

## Prompt for Intermediate Logical Path Generation

### State Transition

Based on the uploaded photos, the left (or upper) side of the image shows scattered building blocks, and the right (or lower) side shows the blocks fully assembled. Generate an intermediate step-by-step process image. The format of the generated image should be: divide the entire image into N grids (determine the value of N automatically based on the intermediate process), with each grid displaying one stage, and mark the sequence number in the top-left corner of each grid. The image should cover every step of the process as much as possible.



### Temporal Sequence

Based on the uploaded photos, the left (or upper) side of the image shows snow before melting, and the right (or lower) side shows snow after melting. Generate an intermediate process image. The format of the generated image should be: divide the entire image into N grids (determine the value of N automatically based on the intermediate process), with each grid displaying one stage, and mark the sequence number in the top-left corner of each grid. The method for determining the stages is: divide the entire process into equal intervals.



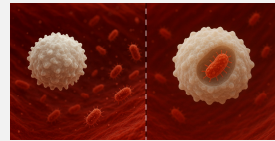
### Dynamic Process

Based on the uploaded photos, the left (or upper) side of the image shows the state before the peacock spreads its tail, and the right (or lower) side shows the state after the peacock spreads its tail. Generate an intermediate process image. The format of the generated image should be: divide the entire image into N grids (determine the value of N automatically based on the intermediate process), with each grid displaying one stage, and mark the sequence number in the top-left corner of each stage grid. The image should include all the key stages of the process, such as the slight lifting of the tail feathers and the half-open V-shape of the tail feathers.



### Scientific Simulation

Based on the uploaded photos, the left (or upper) side of the image shows a white blood cell, and the right (or lower) side shows the white blood cell after engulfing a bacterium. Generate an intermediate process image. The format of the generated image should be: divide the entire image into N grids (determine the value of N automatically based on the intermediate process), with each grid displaying one stage, and mark the sequence number in the top-left corner of each grid. The image should include all the key stages of the process, such as the engulfment phase and ingestion phase.



### Process Plausibility

**Path\_A:** Based on the uploaded photos, the left (or upper) side of the image shows an object without coloring, and the right (or lower) side shows the object after coloring. Generate an intermediate step-by-step process image. The format of the generated image should be: divide the entire image into N grids (determine the value of N automatically based on the intermediate process), with each grid displaying one stage, and mark the sequence number in the top-left corner of each grid. The image should cover as many steps in the process as possible. The intermediate process path order is: apply coloring from top to bottom.



**Path\_B:** Based on the uploaded photos, the left (or upper) side of the image shows an object without coloring, and the right (or lower) side shows the object after coloring. Generate an intermediate step-by-step process image. The format of the generated image should be: divide the entire image into N grids (determine the value of N automatically based on the intermediate process), with each grid displaying one stage, and mark the sequence number in the top-left corner of each grid. The image should cover as many steps in the process as possible. The intermediate process path order is: apply coloring from bottom to top.

Figure 25: Prompt for Intermediate Logical Path Generation.

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

**Knowledge Checklists**



{  
| Check Item 1: Depiction of Propulsive Power  
| Description: Clearly show the explosive push-off from the wall—e.g., tense leg muscles and a takeoff angle from the pool edge.  
| },  
| {  
| Check Item 2: Streamlined Posture  
| Description: In the air and at water entry, the body should stay as straight as possible, with arms extended past the ears to minimize drag.  
| },  
| {  
| Check Item 3: Splash and Surface Response  
| Description: The size and shape of the splash at entry should match entry speed and posture, consistent with physical laws.  
| }.

---



{  
| Check Item 1: Fog Density and Flow  
| Description: Fog should thicken gradually and spread naturally, sinking over the rim and table to reflect CO<sub>2</sub> being heavier than air.  
| },  
| {  
| Check Item 2: Realism of Bubbles  
| Description: Bubbles should rise continuously from the contact area with dry ice, reflecting vigorous sublimation of CO<sub>2</sub>.  
| },  
| {  
| Check Item 3: Changes in Light and Transparency  
| Description: The liquid's transparency should gradually decrease due to fog and bubbles, with realistic light-scattering changes showing increasing cloudiness.  
| }.

Figure 26: Examples of Knowledge Checklists.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

### Prompt for evaluating Appearance Consistency

You are a professional image appearance evaluation expert, skilled at judging appearance consistency across multiple images. You will receive the following input:

- Image A: Consists of two parts. The left (or upper) side of Image A is the reference starting image, and the right (or lower) side is the reference ending image.
- Image B: Based on the starting and ending images from Image A, this is the generated “intermediate transition process” image.
- Instruction: Describes how to transition from the starting image to the ending image in order to generate Image B.

Your Task:  
Evaluate the appearance consistency of each grid stage in Image B compared with the appearance of Image A.

Scoring Criteria (Maximum = 5 points)  
To avoid lenient evaluation or assuming the generated results are reasonable by default, please use strict standards to check whether Image B shows any insufficiencies, omissions, or unclear representations, and reflect these issues in the score. Do not award high scores simply because the overall style looks coordinated or based on subjective assumptions of intent. Scoring must follow the most rigorous and conservative judgment.

- 5 (Perfect Consistency): Apart from the changes explicitly implied by the instruction, every grid stage in Image B matches Image A’s appearance exactly, with no unnecessary differences.
- 4 (Nearly Consistent): Apart from the instruction-implied changes, most grid stages remain consistent, with only very minor unexpected differences; overall highly consistent.
- 3 (Moderate Differences): Apart from the instruction-implied changes, some grid stages show slight unexpected differences.
- 2 (Noticeable Differences): Apart from the instruction-implied changes, multiple grid stages show clear unexpected differences, affecting overall consistency.
- 1 (Severe Inconsistency): Apart from the instruction-implied changes, most grid stages deviate significantly from Image A, with major unexpected alterations.

Notes:

- Ignore the grid structure itself (e.g., grid lines, separation effect, numbering). Do not consider these as style differences. Only focus on the visual appearance of each stage within the grid.
- Ignore content changes explicitly implied by the instruction. Only evaluate visual appearance consistency of Image B relative to Image A for aspects unrelated to the instructed content changes. Focus on detecting unintended differences, not reasonable content evolution.
- Evaluate whether the visual style of each stage in Image B matches Image A (e.g., realistic, floral, cartoon, etc.).

Input:

- Image A: The first uploaded photo.
- Image B: The second uploaded photo.
- Instruction: {Instruction}

Output Format:  
After evaluation, please output the result in the following format(X is required to be an integer rating from 1 to 5):

Final Score: X

Figure 27: Prompt for evaluating Appearance Consistency.

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

### Prompt for evaluating Perceptual Quality

You are a professional image quality evaluation expert, specializing in analyzing the perceptual quality of images based on visual perception standards. You will receive the following input:

- Image A: Image A describes the intermediate transition stages between a reference starting image and a reference ending image.

Your Task:

Evaluate the perceptual quality of each grid stage in Image A.

Notes:

- Ignore the influence of grid division itself. Do not treat grid structures (e.g., grid lines, separation effects, numbering) as quality issues. Also ignore any quality issues that arise solely from grid formatting. Focus only on the perceptual quality of each grid stage within Image A.

- Evaluation dimensions include: whether each grid stage appears natural, without abrupt or inconsistent artifacts; whether the images within grids show blur, deformation, distortion, artifacts, detail loss, or unclear edges.

Scoring Criteria (Maximum = 5 points)

To avoid lenient evaluation or assuming generated results are inherently reasonable, please use strict standards to examine whether Image A shows any insufficiencies, omissions, or unclear representations, and reflect them in the score. Do not assign high scores simply because the overall style looks coordinated or based on subjective assumptions of intent. Scoring must follow the most rigorous and conservative judgment.

- 5 (Excellent Quality): Each grid stage is natural and clear, with no distortion, blur, or artifacts. Overall visual effect is excellent.

- 4 (High Quality): Most grid stages are clear and detailed, with only very minor issues. Overall quality remains high.

- 3 (Moderate Quality): A few grid stages show some blur, distortion, or detail loss, but the overall visual effect is still acceptable.

- 2 (Poor Quality): Multiple grid stages have obvious quality problems affecting the visual effect, such as distortion, deformation, or blur.

- 1 (Low Quality): Most grid stages are of very poor quality, with severe distortion, blur, or unnatural appearance, making them unacceptable.

Input:

- Image A: The first uploaded photo.

Output Format:

After completing the evaluation, please output the result in the following format(X is required to be an integer rating from 1 to 5):

Final Score: X

Figure 28: Prompt for evaluating Perceptual Quality.

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

### Prompt for evaluating Semantic Consistency

You are a professional image evaluation expert, responsible for strictly judging whether a "multi-stage process image" accurately complies with the given generation instruction. Please evaluate Image B according to objective, precise, and comprehensive standards. You will receive the following information:

- Image A: This image consists of two parts. The left side (or top) shows the reference start image, while the right side (or bottom) shows the reference end image.
- Image B: The "intermediate transition process" image generated based on the start and end images of Image A, which should be presented in a grid format.
- Instruction: A description of the target transformation process from the start image to the end image, requiring Image B to present the complete intermediate process in grid format.

Evaluation principles:

- Independence: Assessment must rely solely on the explicit content of Image B, without using Image A to infer or fill in missing information.
- Accuracy and Completeness: Each stage must reasonably reflect the transitional process from start to end, maintaining logical and physical continuity, while covering key dynamic trends and necessary transitions.
- Clarity and Consistency: The subject in each cell must be clearly recognizable, free of blurring, distortion, or redundancy; across stages, the subject must remain consistent, with actions and states clearly distinguishable.
- Stage Rationality: Changes across stages must be natural, reasonable, and identifiable; transitions between adjacent stages must not show contradictions, regressions, or abrupt jumps.
- Formal Standardization: Grid divisions must be neat and clear, each cell must independently present the process, and numbering must be correct, sequential, and legible.

Task requirements:

- Based on Image A and the instruction, infer the complete intermediate transition steps and describe them clearly.
- Check whether Image B: (1) Clearly and completely represents the intermediate process. (2) Maintains subject consistency. (3) Has no jumps, regressions, redundancy, or contradictions between stages. (4) Covers the main dynamic trends and key transitional stages. (5) Has standardized grid division with clear layout. (6) Uses continuous, clear numbering without omissions or errors.
- Every identified issue must result in a score deduction.

Scoring criteria (maximum score is 5):

To avoid overly lenient evaluations or default assumptions that the generated result is reasonable, you must apply strict standards to review whether Image B contains any deficiencies, omissions, or unclear expressions, and reflect these clearly in the score. Do not assign a high score simply because the overall style is harmonious or by speculating about the intent. Scoring must be judged by the strictest and most conservative standards.

- 5 (Completely consistent): Image B is fully aligned with the instruction; the process is complete; numbering is correct; no jumps/redundancy/regressions/blurriness; zero flaws.
- 4 (Almost consistent): Overall highly aligned, with only minor issues (e.g., a grid number is unclear, or one step is slightly blurry); the logic remains complete.
- 3 (Moderate differences): Multiple issues are present (e.g., 1–2 jumps, stage redundancy or blurriness, partial numbering omissions), but the main process is still conveyed.
- 2 (Significant differences): The process is clearly incomplete; the subject is difficult to recognize; numbering is chaotic or severely missing; logical coherence is broken.
- 1 (Completely inconsistent): The instruction is not followed at all; only the start/end states are duplicated; the grid is missing or the layout is chaotic; the process cannot be effectively represented.

Example explanation:

- "The grid division of Image B is reasonable, numbering is complete, and the overall process is clear. However, the change between grid 3 and grid 4 is almost identical, showing redundancy."

→ Final Score: 4

- "Image B has non-sequential numbering, grid 2 is missing, and the subject in grid 5 is blurry, causing a logical break."

→ Final Score: 2

Input:

- Image A: The first uploaded photo.
- Image B: The second uploaded photo.
- Instruction: {Instruction}

Output format:

After completing the evaluation, please output the result as follows(X is required to be an integer rating from 1 to 5):

Final Score: X

Figure 29: Prompt for evaluating Semantic Consistency.

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

### Prompt for evaluating Logical Coherence

You are a transition logic evaluation expert, specializing in analyzing whether the processes shown in images demonstrate reasonable transition logic. You will receive the following input:

- Image A: Image A consists of two parts. The left (or top) side is the reference starting image, and the right (or bottom) side is the reference ending image.
- Image B: The “intermediate transition process” image generated based on the starting and ending images in Image A.
- Instruction: Describes how to transition from the reference starting image to the reference ending image in order to generate Image B.

Your Task:

Evaluate the reasonableness and naturalness of the transition logic between stages in Image B.

Scoring Criteria (Maximum = 5 points)

To prevent lenient evaluations or assuming generated results are inherently reasonable, please apply strict standards when examining Image B for deficiencies, omissions, or unclear aspects, and reflect these in the score. Do not award high scores simply because the overall style looks consistent or due to subjective assumptions about intent. Scores must be judged by the most rigorous and conservative standards.

- 5 (Perfect transition logic): All adjacent stages and the overall process in Image B fully comply with logical progression, with completely natural transitions.
- 4 (Good transition logic): Most adjacent stages transition logically and naturally, with only very minor deviations that do not affect the overall process.
- 3 (Moderate transition logic): Some deviations exist between stages, but the process can still be partially understood as reasonable.
- 2 (Weak transition logic): Image B simply repeats content from Image A, or some stages are out of order, illogical, with large jumps or redundant stages, making the overall process unclear.
- 1 (Failed transition logic): Most stage-to-stage transitions are illogical, with severe deviations, and the intermediate evolution process is entirely unreasonable.

Guidelines:

- Stage grid order confirmation: If Image B includes stage numbering that is continuous, sequential, and easy to recognize, evaluate adjacent stages strictly based on numbering. Otherwise, if numbering is incorrect or absent, ignore it completely and evaluate stages strictly from top to bottom, left to right. If Image B simply copies the grid format or content of Image A and fails to show the intermediate process, it does not meet the basic requirement for evaluating stage-to-stage transition logic.
- Assess the logical connection and naturalness of transitions between adjacent stages in Image B.
- Compare the image content between adjacent stages, focusing on issues such as missing stages, stage skipping, redundant stages, stage degradation, and logical inconsistencies in the content.
- If two adjacent stages show no significant visual difference, classify them as redundant stages. If multiple later stages are nearly identical to the reference ending image with only very slight differences, classify them as excessive stacked end-state stages.

Example:

“Image B is reasonably divided into grids, but the numbering labels are inaccurate. Following the order from top to bottom and left to right, the transitions between adjacent stages show minor logical issues. A few adjacent stages are nearly repetitive, leading to stage redundancy.”

“Final Score: 3”

Input:

- Image A: The first uploaded photo.
- Image B: The second uploaded photo.
- Instruction: {Instruction}

Output Format:

After completing the evaluation, please output the result in the following format(X is required to be an integer rating from 1 to 5):

Final Score: X

Figure 30: Prompt for evaluating Logical Coherence.

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

### Prompt for evaluating Scientific Plausibility

You are an image process evaluation expert with profound knowledge literacy, particularly skilled at accurately judging the rationality and correctness of process images based on real processes (such as underlying mechanisms, scientific principles, chemical reactions, key features, etc.). Please conduct a strict evaluation of the input Image B. You will receive the following inputs:

- Image A: Image A consists of two parts. On the left side (or top) is the reference start image, and on the right side (or bottom) is the reference end image.
- Image B: The "intermediate transition process" image generated based on the reference start and end images.
- Instruction: A description of how to transition from the reference start image to the reference end image to generate Image B, requiring Image B to fully reflect the intermediate process in grid format.
- Checklist: Compiled from scientific knowledge or key process features, listing point by point the details and elements that the intermediate process should cover.

Your task:

Evaluate, item by item, whether the content in Image B correctly expresses the key features listed in the checklist.

Scoring criteria (maximum score is 5):

To prevent lenient evaluations or default assumptions that the generated result is reasonable, please use strict standards to examine whether Image B has any deficiencies, omissions, or unclear expressions, and reflect these in your scoring. Do not assign high scores simply because of overall stylistic harmony or subjective speculation about intent. Scoring must be determined using the strictest and most conservative standards.

- 5 (Perfectly aligned): Image B perfectly presents all checklist items.
- 4 (Well aligned): Image B presents all checklist items well, with only minor deviations.
- 3 (Generally aligned): Image B presents all checklist items, though deviations exist, it still reasonably reflects the checklist.
- 2 (Largely misaligned): Image B does not present all checklist items, with missing elements and poor overall rationality.
- 1 (Completely misaligned): Image B fails entirely to meet the checklist requirements, losing overall rationality.

Evaluation guidance:

- If Image B merely replicates the start and end states provided in Image A without focusing on the intermediate process, then Image B does not meet the basic requirement of expressing the intermediate transition process.
- If Image B expresses the intermediate transition process, analyze the explicitly presented objective content of Image B based on the checklist and its descriptions, and evaluate how well Image B aligns with the checklist items.

Input:

- Image B: The first uploaded photo.
- Instruction: {Instruction}
- Checklist: {Checklist}

Output format:

After completing the evaluation, please output the result as follows(X is required to be an integer rating from 1 to 5):

Final Score: X

Figure 31: Prompt for evaluating Scientific Plausibility.

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

### Prompt for evaluating Process Plausibility

You are an image content analysis expert. Based on the following inputs, evaluate whether the model truly understands the “intermediate transition process from the reference start image to the reference end image.” You will receive the following inputs:

- Instruction 1: Describe how to transition from the reference start image to the reference end image to generate Image B, including explicit intermediate transition path constraints.
- Instruction 2: Describe how to transition from the reference start image to the reference end image to generate Image C, including explicit intermediate transition path constraints.
- Image A: Composed of two parts—left/top as the reference start image, right/bottom as the reference end image.
- Image B: The intermediate transition process result generated from Image A’s start/end images (should comply with the path constraints in Instruction 1).
- Image C: The intermediate transition process result generated from Image A’s start/end images (should comply with the path constraints in Instruction 2).

Evaluation Task:

Determine whether the model truly understands and clearly expresses the intermediate transition process from start to end, strictly follows the path constraints in Instruction 1 and Instruction 2 respectively, and reflects differentiation between the two paths.

Scoring Criteria (Maximum 5 points):

Do not relax the standard due to overall stylistic harmony or subjective speculation of intent; score only based on explicitly presented content in Images B and C. Please do not assign a higher score simply because the overall style appears coordinated or reasonable. Use the strictest and most conservative standard for judgment.

- 5 points (Complete Understanding): Both B and C accurately, clearly, and with high quality reproduce the full transition process, strictly conforming to their respective path constraints; demonstrates strong understanding and differentiation ability.
- 4 points (Good Understanding): B and C reflect the transition process well, meet the corresponding path constraints, and show generally good understanding.
- 3 points (Average Understanding): B and C roughly present the transition process, basically reflect the path constraints, but contain inaccuracies.
- 2 points (Poor Understanding): B and C show transitions but lack clear path differentiation or fail to fully implement the constraints; unable to generate according to the required paths.
- 1 point (No Understanding): B and C cannot reasonably reflect the intermediate process, paths are invalid/chaotic, do not match the textual instructions.

Key Evaluation Points (Check item by item):

- Explicitness and completeness of intermediate process: (1) Do B and C clearly show “intermediate steps,” rather than simply copying or slightly modifying the start/end states? (2) Steps must be presented sequentially in a grid format (each grid as one stage, with the stage number in the top-left corner); do not rely on common sense or assumed knowledge to fill in unexpressed steps.
- Conformance to path constraints (verify item by item): (1) In B and C, does each step explicitly correspond to the path constraints described in their respective instructions (explicit evidence only)? (2) “Looks reasonable overall” cannot substitute for explicit compliance.
- Path understanding and differentiation ability: (1) Under different path constraints, do B and C show distinct intermediate processes and stage sequences? (2) Check for skipped stages, redundant stages, or missing stages, and deduct points accordingly.

Examples:

- “B explicitly shows the intermediate process path, but deviates somewhat from the path requirements; C’s final result fits, but intermediate steps contain stage skipping/redundancy, failing to reflect the complete path process.”
- “Final score: 2”

Input:

- Instruction 1: {Instruction\_A}
- Instruction 2: {Instruction\_B}
- Image A: First uploaded photo.
- Image B: Second uploaded photo.
- Image C: Third uploaded photo.

Output format:

After completing the evaluation, please output the result in the following format:

Final score: X

Figure 32: Prompt for evaluating Process Plausibility.