

---

# Position: Agentic Systems Constitute a Key Component of Next-Generation Intelligent Image Processing

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This position paper argues that the image processing community should broaden  
2 its focus from purely model-centric development to include agentic system design  
3 as an essential complementary paradigm. While deep learning has significantly  
4 advanced capabilities for specific image processing tasks, current approaches face  
5 critical limitations in generalization, adaptability, and real-world problem-solving  
6 flexibility. We propose that developing intelligent agentic systems, capable of  
7 dynamically selecting, combining, and optimizing existing image processing tools,  
8 represents the next evolutionary step for the field. Such systems would emulate  
9 human experts' ability to strategically orchestrate different tools to solve complex  
10 problems, overcoming the brittleness of monolithic models. The paper analyzes key  
11 limitations of model-centric paradigms, establishes design principles for agentic  
12 image processing systems, and outlines different capability levels for such agents.

## 13 1 Introduction

14 Image processing is a longstanding re-  
15 search area in computer vision. We  
16 have a wide variety of image pro-  
17 cessing and editing needs, ranging  
18 from post-photography editing, im-  
19 age restoration, enhancement, to style  
20 transfer. These tasks are inherently  
21 complex due to both the intricate na-  
22 ture of images and the unique aes-  
23 thetic standards and nuanced expec-  
24 tations that humans hold. For a long  
25 time, image processing has been a spe-  
26 cialized technical field managed by  
27 dedicated technicians and artists. Efforts in computer vision have long aimed to provide high-quality  
28 tools that enhance the efficiency and effectiveness of image processing tasks. *The research community*  
29 *strives to develop intelligent, adaptable software that maximizes convenience for users at all levels*  
30 *and fulfills a wide range of image processing needs.*

31 Early image processing algorithms were typically designed for specific types of problems, making  
32 them part of a broader pipeline or a standalone tool [5]. Professionals often need to configure and  
33 combine multiple processing steps to address particular image processing challenges. In the past  
34 decade, deep learning has driven a major leap in image processing, significantly improving the  
35 quality of individual tasks while introducing a more generalized, intelligent paradigm [14, 65]. The

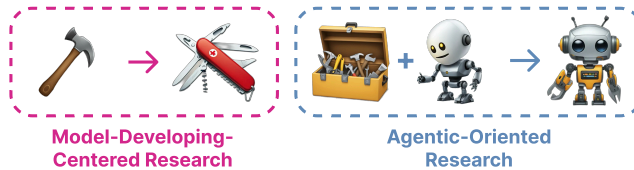


Figure 1: The existing research paradigm focuses on developing more powerful and multi-functional image processing models. In contrast, we advocate a new research paradigm centered on building agentic systems. Our goal is to create an agent that can integrate and leverage these models to achieve higher levels of intelligence, automation, and generality.

36 industry has gradually shifted from constructing image processing pipelines to training end-to-end  
37 deep learning models that replace complex pipelines [3]. These deep learning models are now used  
38 not only to solve isolated problems but also to establish a general, multi-task, and intelligent solution  
39 that can operate effectively in diverse, real-world conditions [31, 50, 62, 68]. The advent of deep  
40 learning and artificial intelligence has made the vision of a general, intelligent, software-based “image  
41 processing assistant” seem closer than ever, though it remains just out of reach. **Significant research  
42 efforts have focused on the current paradigm, which is predominantly centered on developing  
43 various deep image processing models, as shown in Figure 1 (left).**

44 Nevertheless, the limitations of deep networks are emerging, and continuing within the existing  
45 research paradigm makes overcoming these constraints challenging. Firstly, these models face issues  
46 with generalization, as they perform well on test data similar to their training data but struggle on  
47 test data that deviates significantly from it [20]. Secondly, deep models capable of handling a wide  
48 range of degradation scenarios often compromise on quality and generalization [66]. Those that  
49 excel in specific degradations may lack generalizability, while models that handle a broad spectrum  
50 of tasks may not deliver peak performance on any single task [31, 67]. These challenges suggest  
51 that relying on a single model or fixed process for image processing may not effectively address the  
52 dynamic and complex real-world problems. Interestingly, despite the limitations of current image  
53 processing models, human artists and image editing professionals can still leverage these models and  
54 tools – often very simple ones, like basic operations in PhotoShop – to accomplish complex tasks that  
55 even the most advanced models cannot achieve. Emulating the dynamic and adaptive ways in which  
56 humans use these image processing tools could be a crucial step toward making image processing  
57 more intelligent and general. After all, no matter how powerful or multi-functional a tool may be, it  
58 still requires a capable operator – human or otherwise – to realize its true potential, see Figure 1 right.

59 **In this position paper, we advocate for a new research paradigm centered on agentic-oriented  
60 image processing systems, offering a more autonomous, adaptable, and intelligent alternative  
61 to current methods.** We begin by discussing the core capabilities required for intelligent image  
62 processing systems and the challenges that the current paradigm faces in Sec. 2. We then introduce  
63 the concept of AI Agents in Sec. 3, exploring their fundamental principles and role in intelligent  
64 systems. In Sec. 4, we extend this discussion to agentic image processing systems, analyzing how  
65 existing methods can incorporate varying degrees of agentic features to enhance generality and  
66 intelligence. We also examine the key characteristics of agentic image processing systems and outline  
67 different levels of agentic capability. Recognizing that large language models have become pivotal in  
68 the study of intelligent agents, we also explore how language models and multi-modal techniques  
69 may shape the future of image processing. Moreover, in Sec. 5, we highlight that there remains  
70 further room for exploration in certain critical attributes that determine a system’s level of intelligence  
71 and generalization.

72 **Alternative Views.** The prevailing view holds that continued progress in developing models  
73 – through scale and improved architectures – could eventually overcome all the limitations and  
74 subsume the proposed agentic capabilities. While these views have merit, they underestimate the  
75 fundamental mismatch between static model architectures and the dynamic, compositional nature  
76 of real-world image processing requirements. Hybrid approaches combining foundation models  
77 with agentic components may offer a viable middle ground, but system intelligence requires explicit  
78 architectural support beyond current paradigms.

## 79 **2 Backgrounds**

### 80 **2.1 Intelligent Image Processing**

81 The core of intelligent image processing lies in developing intelligent and efficient algorithms that  
82 enable computers to automatically and accurately process images in various conditions to meet visual,  
83 psychological, or other needs. Ultimately, its goal is to build a “software employee” capable of  
84 automatically, intelligently, and effectively completing various image processing tasks. This is a  
85 highly visionary goal, and the community has long been approaching it from different angles in  
86 an attempt to simplify this challenging problem. Initially, image processing methods were mainly  
87 extensions of signal processing techniques applied to two-dimensional image signals, focusing on  
88 specific image processing operations. Entering the 21st century, with advancements in computing  
89 ability, many methods based on image priors and optimization have emerged but remain limited to

90 specific task scenarios. In the past years, the rise of machine learning and deep learning [32, 14]  
91 has propelled significant progress in the field of image processing. Particularly, the introduction  
92 of neural networks has led to breakthrough results in various image processing tasks. Data-driven  
93 methods not only allow for multiple image tasks to be handled within the same algorithmic framework  
94 but also make multitask integration and general-purpose image processing possible. Researchers  
95 have once again embraced the vision of intelligent image processing, and constructing a “software  
96 employee” capable of handling all tasks by unifying image processing tasks seems to have become a  
97 feasible direction. This position paper focuses on the core challenges of realizing this vision and the  
98 approaches to overcome them.

Specifically, an intelligent image processing system, a “software employee”, needs to possess at least the following core capabilities:

- *Generality*: The system should be able to handle a wide range of diverse tasks without requiring separate models for each one, nor relying on extensive domain-specific training data or explicit task-specific instructions.
- *Autonomous*: The system should minimize reliance on user operations and supervision. It demonstrates proactivity by leveraging prior experience to explore new strategies without requiring explicit instructions.
- *Intelligence*: The system can adaptively adjust its processing strategies based on the semantic content and quality of the input, and user instructions. The system should demonstrate its intelligence and complexity in the processing workflow or in the final outcomes.
- *User Interaction and Feedback*: The system should facilitate clear, continuous, and user-friendly communication.
- *Self-Evolution and Creativity*: The system should generate meaningful and innovative solutions, going beyond straightforward problem-solving to provide novel approaches, insights, or outputs that showcase originality. Additionally, the system can continuously evolve and learn from new data, experiences, and user interactions.

99

## 100 2.2 Why Intelligent Image Processing is Challenging?

101 However, the current mainstream research paradigm, which centers around the development of deep  
102 learning models, struggles to align with the aforementioned vision.

103 **The Challenge of Achieving Generality.** Unlike high-level image understanding tasks, image  
104 processing tasks have both input and output as images that require precise pixel-level correspondence.  
105 The information needed for image processing tasks is not as specific as in image understanding. In  
106 image understanding, the model abstracts the image, extracts main features, and aligns them with  
107 semantics expressible in human terms. Although we hope that advanced deep networks for image  
108 processing can also learn the “semantics” of images, it is challenging to accurately describe local  
109 image details semantically. In fact, image processing networks do not learn semantics [19, 37, 22]  
110 but instead learn certain image transformations and overfit to training degradations [20, 39]. This is  
111 determined by the training paradigm of deep image processing models. Therefore, essentially, current  
112 image processing networks are not intelligent.

113 This leads to the next issue: the differences between various image processing tasks are also distinct  
114 from other types of tasks. Generally speaking, a specific image transformation or degradation can  
115 define an entirely new task. The differences between image transformations or degradations can be  
116 very subtle, and they can also be compounded to create almost unlimited types of transformations or  
117 degradations, resulting in virtually infinite image processing tasks. Due to deep models overfitting  
118 to the training set [8], tasks beyond the training scope cannot be well addressed [38]. This greatly  
119 limits the ability of current image processing methods to solve general problems. Worse yet, because  
120 collecting training images in the real world is extremely difficult, most research can only train on  
121 synthetic data, which further leads to generalization issues in practical applications.

122 Some methods attempt to include as many tasks as possible in the training set and train a sufficiently  
123 large model to achieve generality for common tasks [31, 11], even hoping that increasing the number  
124 of tasks will enable the network to generalize. However, these models have been proven to have a  
125 trade-off between the range of tasks and processing performance [66]. It’s challenging to expand

126 the task range while keeping image processing performance from significantly declining. All these  
127 issues make constructing image processing systems with general capabilities highly challenging.

128 **The Challenge of Developing Intelligence.** Beyond the requirements of generality, we are increas-  
129 ingly emphasizing the intelligence these image processing systems exhibit. Firstly, we hope that  
130 image processing systems can explicitly perceive image content and perform targeted processing  
131 based on that content. For example, generating corresponding fur on animals or inferring and com-  
132 pleting blurry or missing objects. Existing research indicates that end-to-end supervised deep image  
133 processing models do not possess this characteristic [37], but methods based on pre-trained generative  
134 models have demonstrated related capabilities and have thus achieved good results [62]. Secondly,  
135 we expect intelligent image processing systems to adaptively adjust processing strategies based on  
136 different input types or qualities, and even have the ability to make complex decisions based on  
137 specific image content. For instance, the system can automatically select the optimal denoising,  
138 enhancement, or restoration methods according to the image’s resolution, lighting conditions, or  
139 noise levels. Additionally, we hope that image processing systems can dynamically understand users’  
140 complex needs. Currently, users need to select tools and set parameters based on their own expertise  
141 before obtaining results; this process does not reflect the system’s intelligence. An intelligent system  
142 can accept user feedback or instructions to make dynamic adjustments in subsequent processing.  
143 These requirements have been mentioned to varying degrees in image processing research, but none  
144 have been explored in depth.

145 **The Challenge of Balancing Autonomy, User Interaction, and Creativity.** Existing approaches  
146 often fall into two extremes. On one hand, fully autonomous methods – such as end-to-end models  
147 – can quickly complete tasks but tend to overlook subtle user preferences, resulting in a rigid,  
148 one-size-fits-all automation. Automatic denoising may eliminate intentionally added artistic grain,  
149 and style transfer algorithms can homogenize diverse creative visions. Given the broad range and  
150 complex demands of image processing tasks, achieving consistently high-quality results proves  
151 challenging with these models. The end result is that people still need to pick and combine the results  
152 of different models, thus losing this automaticity. On the other hand, heavily manual interfaces  
153 impose a significant technical burden on users. Professional software like Photoshop requires  
154 extensive manual intervention and expert knowledge, which conflicts with the goals of ease of use  
155 and accessibility. Moreover, many existing approaches rely on single-model solutions with limited  
156 interactivity; more semantic, higher-level, and varied interaction mechanisms are needed to facilitate  
157 seamless communication between the user and the system.

158 Furthermore, existing methods also struggle to foster genuine creativity. Here, “creativity” goes  
159 beyond generating novel content via generative models [62]; it also involves discovering innovative  
160 ways to repurpose existing tools and deepening our understanding of them. As image processing  
161 evolves from mere technical correction into a creative medium, bridging this gap demands systems  
162 that not only “see” the pixels but also interpret the cultural, emotional, and contextual layers – a  
163 frontier that remains largely unexplored in current technology.

### 164 3 What is AI Agent?

165 An agent is a program designed to achieve its goals by perceiving the environment and interacting  
166 with it through available tools. These agents can operate autonomously without human intervention  
167 and proactively work towards their objectives [17, 34, 56]. From a design perspective, agent-  
168 based systems naturally fulfill our demand for automation. The various image processing models  
169 we develop can be regarded as tools of different scales and purposes, while the agent acts as  
170 the “coordinator” that actively orchestrates these tools, as shown in Figure 1 left. Early agent  
171 programs largely relied on symbolic methods [17] and reinforcement learning [24, 46, 63]. In  
172 recent years, however, agent systems powered by Large Language Models (LLMs) have achieved  
173 transformative progress [58, 7]. By training on massive text corpora through next-token prediction,  
174 LLMs demonstrate powerful knowledge transfer and logical reasoning abilities [4, 41, 1, 47, 25, 15],  
175 showcasing considerable potential in complex reasoning [51, 30], step-by-step planning [57, 58], and  
176 domain-specific knowledge applications [40, 21]. Compared to traditional reinforcement learning  
177 agents, LLM-based agents maintain long-term planning and simultaneously leverage broad general  
178 knowledge, thereby exhibiting more human-like cognitive characteristics [1]. From a cognitive  
179 standpoint, the fundamental responses of an LLM can be likened to “System 1,” characterized by  
180 rapid, automatic thinking, whereas more advanced composite agent systems emulate “System 2,”

181 which involves deliberate, reflective reasoning [57, 35, 28]. Recent research has explored diverse  
182 agent architectures that enhance LLM-based problem-solving through structured mechanisms [43, 49]  
183 – such as tree- or graph-based search strategies [2, 57], external tool integration [44, 52], memory  
184 retrieval systems [70, 42], and error-driven learning processes [45, 58]. By combining an LLM’s  
185 reasoning capabilities with structured problem-solving frameworks, these approaches show strong  
186 potential for tackling complex tasks [16].

187 Notably, pioneering efforts have employed LLM agents across various domains, all striving to create  
188 automated systems capable of proactively tackling a broad spectrum of challenges, aligning with  
189 the vision outlined in this position paper. For instance, frameworks like HuggingGPT [44] and  
190 Visual ChatGPT [44] leverage LLMs as multi-modal task controllers, integrating them with model  
191 libraries to decompose and solve diverse tasks; frameworks like OctoPack integrate LLMs with  
192 specialized toolsets, achieving significant performance gains in fields like medical image processing  
193 [40]. Advancements have also highlighted the effectiveness of LLM agents in tackling complex image  
194 processing tasks, achieving remarkable results [69, 9]. These advancements collectively highlight the  
195 transformative potential of LLM-based agents in addressing complex multi-modal challenges.

## 196 4 Agentic Image Processing System

197 The initial step toward agentic image processing involves acknowledging the fundamental reality that,  
198 **regardless of how advanced your image processing model is, carefully chosen preprocessing,**  
199 **postprocessing, or application-specific techniques/tricks can often enhance its performance.**  
200 For instance, certain severe degradations cannot be fully restored by a single pass through an image  
201 restoration model; applying the model iteratively to its own outputs can yield further improvements.  
202 Additionally, some degradations may lie beyond the training scope of the model, and introducing  
203 deliberate additional blurring before restoration can significantly mitigate these challenging cases.  
204 There exist numerous possibilities for such operations, and in practical applications, users frequently  
205 leverage these techniques to maximize performance.

### 206 4.1 Paradigms of Current Image Processing System

207 While this paper is the first to advocate for the construction of an agentic system to address challenges  
208 in intelligent image processing, traces of agentic thinking have already emerged, to varying degrees,  
209 in previous studies. We begin by examining the embodiment of agentic concepts behind the design of  
210 existing methods, adopting a perspective that progresses from simple to complex. Figure 3 provides a  
211 schematic illustration of these paradigms, offering a visual aid for better understanding.

212 **End-to-End** models are the most common paradigm in image processing research. Given an input, an  
213 end-to-end model produces a corresponding output. This category encompasses optimization-based,  
214 filter-based, and deep network models, with a focus on end-to-end deep network models for intelligent  
215 models. The standard approach involves collecting images that need processing along with their  
216 corresponding target images to form training image pairs, and then training the image processing  
217 model on this basis. This paradigm is the least agentic, and due to the following reasons, it has  
218 limitations in terms of generality and intelligence: Due to the limitations discussed in Sec. 2, no  
219 single model can simultaneously achieve both broad image processing capabilities and outstanding  
220 results. If a model is designed to be sufficiently “general,” it will inevitably come at the cost of  
221 reduced performance on specific tasks.

222 **Pipeline** paradigm typically decomposes complex and difficult-to-model-at-once image processing  
223 problems into multiple independent processing steps. The main advantage of this approach is that it  
224 can effectively break down complex tasks into more manageable subtasks, allowing for the creation  
225 of new tasks through the combination of a limited number of image processing/operations [5, 3].  
226 The modular design also equips the pipeline paradigm with high flexibility and scalability, enabling  
227 the system to be adjusted and updated according to specific needs. This makes it convenient to  
228 integrate new technologies or algorithms into the existing framework. For example, users can directly  
229 replace the denoising step with the latest denoising algorithm without redesigning the entire pipeline.  
230 Pipeline design is a typical idea of people to solve complex problems by combining simple tools.

231 Although pipeline models have advantages in handling complex tasks, their agentic level is still  
232 relatively low because each step is pre-defined based on practical applications and is difficult to adjust

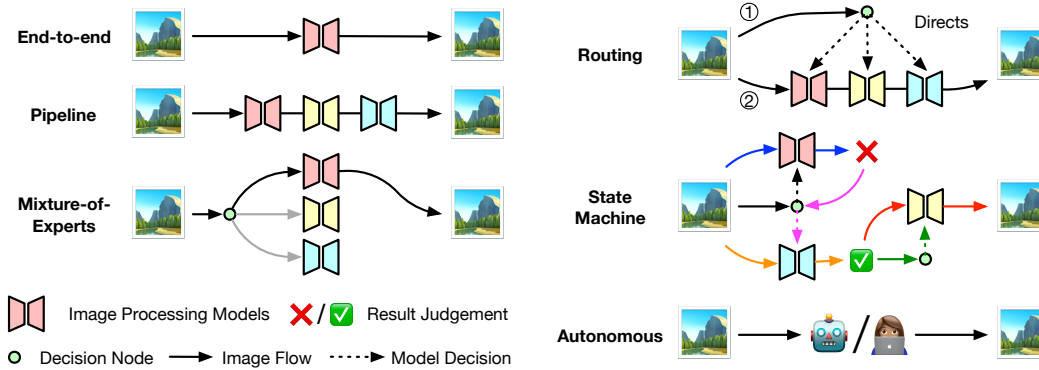


Figure 2: How image processing systems can embody different levels of agentic to enhance their generality and intelligence.

233 according to the diversity of inputs. Moreover, pipeline models often rely heavily on manual design,  
 234 as the task decomposition and execution order within the pipeline significantly impact the results  
 235 [63, 9, 69]. Therefore, they are used for specialized solutions to specific real-world problems rather  
 236 than aiming for pursuing generality or higher levels of intelligence. However, by integrating different  
 237 steps, the pipeline approach expands the application boundaries of image processing algorithms. This  
 238 characteristic is a core advantage that agentic image processing systems can leverage.

239 **Mixture-of-Experts (MoE)** is another paradigm that broadens the task range of image processing  
 240 systems and enhances performance by integrating the capabilities of multiple models. A single  
 241 model is constrained by the trade-off between task coverage and processing effectiveness, making it  
 242 difficult to efficiently handle each individual task while covering a large number of tasks. Similar  
 243 phenomena have been observed in other large-scale model practices [6, 13, 12, 26]. To overcome  
 244 these limitations, MoE typically introduces multiple expert models, each focusing on a specific  
 245 task. The system dynamically selects the most suitable expert based on the input data or the task  
 246 requirements, or combines the outputs of multiple expert models to optimize processing performance.  
 247 This approach not only achieves a balance between task coverage and processing effectiveness but  
 248 also allows for flexible adaptation to new tasks or improvement of performance on specific tasks by  
 249 adjusting and replacing expert models. Therefore, MoE becomes an effective means to achieve task  
 250 breadth while ensuring processing depth. Although MoE achieves a certain degree of agentic through  
 251 proactive model selection, its generality and intelligence still depend on the performance of each  
 252 expert model within the system. Since we cannot infinitely expand the number of expert systems, and  
 253 there still exist problems that individual models cannot effectively solve, the generality of the MoE  
 254 paradigm remains quite limited.

255 **Routing** is a combination of the Pipeline and MoE paradigms, potentially integrating the advantages  
 256 of both. Similar to the MoE paradigm, the routing paradigm selects corresponding processing paths  
 257 for input images to achieve targeted processing. However, unlike MoE, the routing paradigm selects  
 258 a Pipeline composed of multiple models to maximally expand the range of feasible tasks. In essence,  
 259 the routing paradigm automatically devises dedicated pipelines for different input image tasks and  
 260 invokes the corresponding models. In other words, routing makes a “plan” for each input and executes  
 261 it [23, 63]. The routing paradigm further enhances the system’s agentic; when the decision-making  
 262 methods are sufficiently accurate and robust, this paradigm can greatly expand the potential task  
 263 coverage, thereby improving its generality. Since the decision-making process requires a deeper  
 264 understanding of the images, the routing paradigm also possesses higher intelligence. However, once  
 265 the path is determined, the outcome is already fixed. If an issue arises in an intermediate step, the  
 266 routing approach cannot backtrack to address the problem at that specific step.

267 **State Machine.** Building on the foundation of the routing paradigm, state machines further expand to  
 268 allow more fine-grained control over the processing flow. Similar to routing, a state machine produces  
 269 a complex execution plan to conduct multiple image processing steps. However, due to the complexity  
 270 of images and the variety of image processing operations, it is often not feasible to directly determine  
 271 an optimal plan or parameter set in one run. In contrast to routing, the most notable feature of a  
 272 state machine is its intelligent flow control: the system can reason and autonomously decide whether  
 273 to proceed to the next step, adopt the current result or plan, or even undo the previous operation.

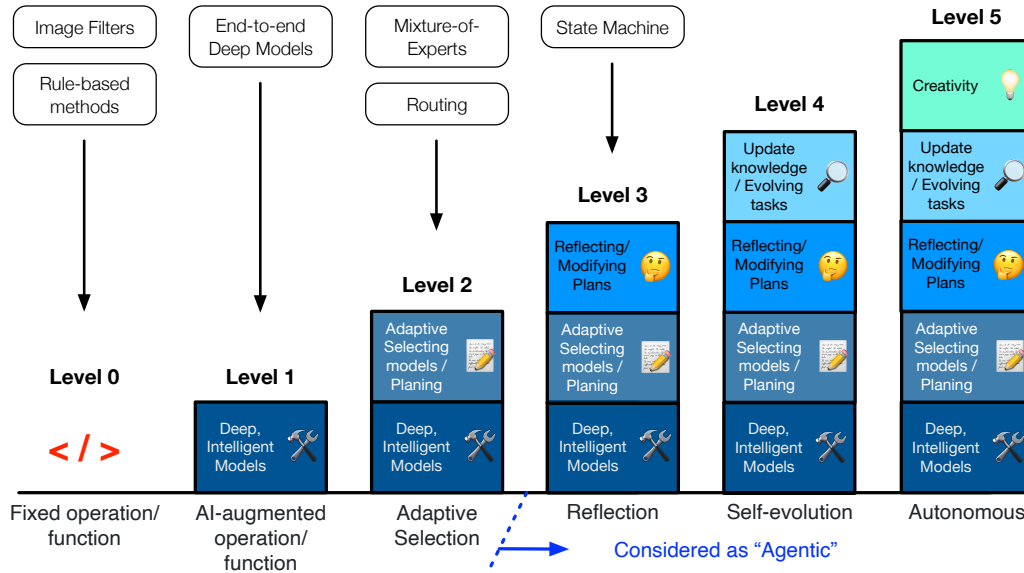


Figure 3: Levels of agentic capability in image processing systems, illustrating the progression from fixed rule-based methods (Level 0) to fully autonomous and creative systems (Level 5). Each level builds on the previous one, adding layers of adaptability, reflection, self-evolution, and creativity.

274 Essentially, the process by which humans solve problems can also be viewed as a highly flexible  
 275 state machine. Some pioneering studies have already adopted this paradigm to build intelligent agent  
 276 systems, demonstrating their remarkable intelligence and potential [69, 9, 36].

## 277 4.2 Characteristics of an Agentic System

278 Based on the above analysis, several core design principles of agentic thinking are already present in  
 279 prior works to varying degrees. To further advance this approach, we summarize below the potential  
 280 features of an agentic system – features that can significantly improve the system’s intelligence,  
 281 generality, and ease of use:

- 282 • *Proactive and Autonomous Problem-Solving*: An agentic system autonomously senses challenges,  
 283 explores different models and methods, and dynamically adjusts strategies without relying on  
 284 further human instructions. This allows for flexible and efficient image processing even in complex  
 285 scenarios.
- 286 • *Integration of Multiple Models/Tools*: Rather than depending on a single model for complex image  
 287 processing tasks, an agentic system can strategically combine multiple models or tools according  
 288 to the specific task requirements or image characteristics. Even “all-in-one” large models can be  
 289 combined with other operations or models to broaden coverage and improve performance.
- 290 • *Adaptive, Context-Aware Strategies*: An agentic system tailors its processing strategy based on the  
 291 specific content or characteristics of the input image instead of applying a fixed pipeline. In other  
 292 words, the system reasons about the image to make informed decisions.
- 293 • *Modular Architecture*: An agentic system often consist of multiple functional modules that work  
 294 together – commonly including Perception, Reasoning, Action, and Reflection. Data and results  
 295 flow through these modules, forming a coherent and synergistic workflow.
- 296 • *Easy Extensibility*: An agentic system should be readily extensible, allowing new features, tools  
 297 or modules to be added without large-scale retraining or constructing massive new datasets. This  
 298 flexibility enables the system to adapt to evolving requirements more effectively.
- 299 • *Continuous Reflection and Improvement*: An agentic system incorporates reflection mechanisms  
 300 to evaluate and refine their performance. This iterative learning process ensures the system can  
 301 leverage real-world usage data for sustained improvement over time.

### 302 4.3 Levels of Agentic in Image Processing

303 It is clear that “agentic” is not a binary concept but rather exists on a continuum. Drawing inspiration  
304 from the levels of autonomous driving [29], we propose a reference framework for classifying the  
305 agentic levels of image processing systems into six tiers, as shown in Figure 2. These levels reflect  
306 different characteristics that such systems may exhibit:

- 307 • **Level 0 (Fixed operation/function).** Methods at this level only provide basic, fixed trans-  
308 formations and processes. Regardless of the input image, they perform operations strictly  
309 according to predefined rules, such as filter-based or rule-based transformations. This stage  
310 exhibits almost no intelligence or generality.
- 311 • **Level 1 (AI-augmented operation/function).** While still focused on specific tasks, systems  
312 at this level go beyond simple rules by incorporating complex patterns learned through deep  
313 learning or other data-driven approaches. Although their performance surpasses Level 0, they  
314 remain limited in generalization.
- 315 • **Level 2 (Adaptive Selection).** Starting from this level, image processing systems no longer  
316 rely on a single model or tool. Instead, they can adaptively select and integrate different  
317 models, thus expanding the range of tasks they can handle. The ability to choose different  
318 processing strategies based on the input image demonstrates a certain degree of intelligence.
- 319 • **Level 3 (Reflection).** This level introduces more freedom in workflow control and the ability  
320 to reflect on output results, allowing systems to flexibly adjust strategies and processes.  
321 Through reflection and iterative adjustments, systems at this level can tackle a wider variety of  
322 problems.
- 323 • **Level 4 (Self-evolution).** At this level, agentic surpasses what fixed architectures can achieve.  
324 Systems can continuously learn and evolve from large amounts of data and experience,  
325 distilling new knowledge to solve previously unsolvable problems. They may even modify or  
326 advance their own workflows.
- 327 • **Level 5 (Fully Autonomous).** At the highest level, agents can execute image processing  
328 tasks autonomously without user intervention. In addition to incorporating all capabilities of  
329 the previous levels, they possess a degree of creativity, enabling innovative problem-solving  
330 approaches (e.g. discovering new tricks that the people who created the agent don’t know  
331 about.). As a result, they can potentially replace human experts and approach a form of  
332 artificial general intelligence (AGI).

### 333 4.4 The Role of LLMs

334 Incorporating agentic design enables greater autonomy and adaptability. Current image processing  
335 paradigms can partially fulfill these objectives if designed with sufficient complexity (e.g., approaches  
336 based on reinforcement learning [23] or expert systems [64]). However, due to the inherent complexity  
337 of image processing tasks and the ambiguity of their descriptions, existing paradigms struggle to  
338 further develop and leverage agentic capabilities. LLMs have demonstrated strong adaptability when  
339 dealing with open-domain problems.

340 When an agent autonomously tackles complex tasks, it often faces numerous scenarios arising  
341 from the interplay of diverse factors. Given the complexity of image processing systems, these  
342 scenarios cannot be easily summarized or handled with simple rules. Traditional methods rely on  
343 predefined rules and features, which become insufficient in the face of combinatorial explosions.  
344 In contrast, LLMs can perform language-based reasoning and planning for each situation, offering  
345 remarkable generalization capabilities. Through this reasoning mechanism, abstract and unstructured  
346 demands can be mapped to specific image processing models or tools and translated into executable  
347 steps. Moreover, the workflow can be dynamically adjusted based on real-time feedback – for  
348 instance, rolling back or modifying the previous step. LLMs may even derive innovative new model  
349 combinations.

350 In a multi-modal setting, LLMs further provide the system with an “intelligent eye,” enabling it to  
351 extract semantic information at multiple levels from abstract visual signals, far beyond what non-LLM  
352 approaches can achieve [61, 60, 53, 59]. Finally, natural language dialogue has proven to be an  
353 efficient and user-friendly channel for interaction. By employing LLMs, the system can engage in  
354 more flexible conversations with users, offer feedback, and accept instructions, thereby significantly  
355 enhancing both usability and extensibility.

## 356 5 Core Problems Demanding Further Study

357 Building intelligent agentic systems is a novel direction with many core challenges that require  
358 in-depth exploration. We analyze key issues that warrant attention in future research. Due to space  
359 constraints, we focus on two potential directions here and discuss additional topics in the appendix.

### 360 5.1 Cognitive Architecture of Image Processing

361 The foundation for building more complex agentic  
362 systems lies in designing their overall **cognitive ar-**  
363 **chitecture**<sup>1</sup>. A cognitive architecture refers to how  
364 the system “thinks” – in other words, the flow of code,  
365 prompts, and LLM calls that accept user input and exe-  
366 cute operations or generate responses. Designing a  
367 cognitive architecture involves contemplating the ab-  
368 stract processes by which an intelligent agent solves  
369 problems at different levels. It’s the methodology an  
370 agent uses to address a certain class of problems.

371 To facilitate understanding, we can start with the  
372 abstracted process of humans performing image pro-  
373 cessing or PhotoShop editing tasks. Figure 4 illus-  
374 trates an example of a personified cognitive archite-  
375 cture for an image processing system, which is also  
376 the architecture used by Zhu et al. [69] and Chen et  
377 al. [9]. In this architecture, the interaction between the system and tools is abstracted into five stages:  
378 **Perception, Scheduling, Execution, Reflection, and Rescheduling.** Specifically, the Perception  
379 stage acts as the agent’s “eyes,” extracting necessary information from the input image. The Schedul-  
380 ing stage functions like the “brain,” making judgments and formulating plans based on the acquired  
381 information and existing knowledge. The Execution stage represents “action,” carrying out specific  
382 operations according to the plan. The Reflection stage evaluates whether the intermediate results  
383 meet expectations. If they do, the agent proceeds with subsequent plans; if not, the Rescheduling  
384 stage considers the failed results and modifies the original plan.

385 However, this intuitive cognitive architecture still leaves much room for improvement in the archi-  
386 tectural research of image processing agent systems. For instance, for more specific problems, how  
387 should we design their cognitive architectures to meet the need for more refined control? For tasks  
388 requiring higher generality, how can we abstract a sufficiently general process to encompass a wider  
389 range of possible tasks? How can we systematically explore, distill, and abstract human problem-  
390 solving strategies into foundational principles for designing cognitive architectures? Furthermore,  
391 how can we create a cognitive architecture that surpasses human limitations, optimized specifically  
392 for intelligent image processing agent systems?

### 393 5.2 More Problems Demanding Further Study

394 Due to space limitations, we discuss four additional topics in the appendix: *Image Quality Assessment*  
395 *and Content Analysis*, *Knowledge Acquisition and Infusion*, *Human-Computer Interaction in Agentic*  
396 *Systems*, and *Exploitative Learning, Self-Evolution & Creativity*.

## 397 6 Conclusion

398 The evolution from task-specific models to agentic image processing systems marks a fundamental  
399 shift in addressing real-world complexity through dynamic tool orchestration rather than monolithic  
400 architectures. By embedding human-like adaptive reasoning into operational frameworks, such  
401 systems transcend current generalization limits while preserving specialized model strengths.

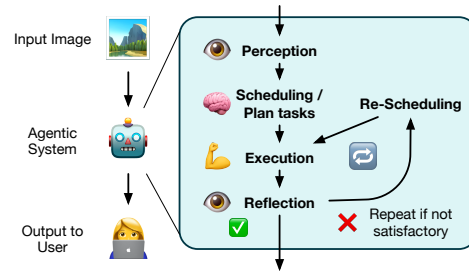


Figure 4: Cognitive architecture for image processing systems, illustrating the iterative process of perception, scheduling, execution, reflection, and rescheduling to achieve satisfactory results.

<sup>1</sup>The term “cognitive architecture” has a rich history in neuroscience and computational cognitive science. It refers both to theories about the structure of human thought and to computational implementations of these theories. Here, we borrow this concept but do not specifically refer to its original meaning.

## 402 References

- 403 [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
404 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv*  
405 *preprint arXiv:2303.08774*, 2023.
- 406 [2] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi,  
407 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving  
408 elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial*  
409 *Intelligence*, volume 38, pages 17682–17690, 2024.
- 410 [3] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unpro-  
411 cessing images for learned raw denoising. In *Proceedings of the IEEE/CVF conference on computer vision*  
412 *and pattern recognition*, pages 11036–11045, 2019.
- 413 [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
414 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.  
415 *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 416 [5] Mark Buckler, Suren Jayasuriya, and Adrian Sampson. Reconfiguring the imaging pipeline for computer  
417 vision. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 975–984, 2017.
- 418 [6] Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. A survey on mixture of  
419 experts. *arXiv preprint arXiv:2407.06204*, 2024.
- 420 [7] Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F Karlsson, Jie Fu, and Yemin Shi.  
421 Autoagents: A framework for automatic agent generation. *arXiv preprint arXiv:2309.17288*, 2023.
- 422 [8] Haoyu Chen, Jinjin Gu, Yihao Liu, Salma Abdel Magid, Chao Dong, Qiong Wang, Hanspeter Pfister, and  
423 Lei Zhu. Masked image training for generalizable deep image denoising. In *Proceedings of the IEEE/CVF*  
424 *Conference on Computer Vision and Pattern Recognition*, pages 1692–1703, 2023.
- 425 [9] Haoyu Chen, Wenbo Li, Jinjin Gu, Jingjing Ren, Sixiang Chen, Tian Ye, Renjing Pei, Kaiwen Zhou,  
426 Fenglong Song, and Lei Zhu. Restoreagent: Autonomous image restoration agent via multimodal large  
427 language models. *Advances in Neural Information Processing Systems*, 2024.
- 428 [10] Haoyu Chen, Wenbo Li, Jinjin Gu, Jingjing Ren, Haoze Sun, Xueyi Zou, Zhensong Zhang, Youliang Yan,  
429 and Lei Zhu. Low-res leads the way: Improving generalization for super-resolution by self-supervised  
430 learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
431 25857–25867, 2024.
- 432 [11] Xiangyu Chen, Yihao Liu, Yuandong Pu, Wenlong Zhang, Jiantao Zhou, Yu Qiao, and Chao Dong.  
433 Learning a low-level vision generalist via visual task prompt. *arXiv preprint arXiv:2408.08601*, 2024.
- 434 [12] Zitian Chen, Yikang Shen, Mingyu Ding, Zhenfang Chen, Hengshuang Zhao, Erik G Learned-Miller, and  
435 Chuang Gan. Mod-squad: Designing mixtures of experts as modular multi-task learners. In *Proceedings*  
436 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11828–11837, 2023.
- 437 [13] Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng,  
438 Xingkai Yu, Y Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts  
439 language models. *arXiv preprint arXiv:2401.06066*, 2024.
- 440 [14] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep  
441 convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307,  
442 2015.
- 443 [15] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,  
444 Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint*  
445 *arXiv:2407.21783*, 2024.
- 446 [16] Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori,  
447 Yusuke Noda, Demetri Terzopoulos, Yejin Choi, et al. Agent ai: Surveying the horizons of multimodal  
448 interaction. *arXiv preprint arXiv:2401.03568*, 2024.
- 449 [17] Stan Franklin and Art Graesser. Is it an agent, or just a program?: A taxonomy for autonomous agents. In  
450 *International workshop on agent theories, architectures, and languages*, pages 21–35. Springer, 1996.
- 451 [18] Jinjin Gu, Haoming Cai, Haoyu Chen, Xiaoxing Ye, Jimmy Ren, and Chao Dong. Image quality assessment  
452 for perceptual image restoration: A new dataset, benchmark and metric. *arXiv preprint arXiv:2011.15002*,  
453 2020.

- 454 [19] Jinjin Gu and Chao Dong. Interpreting super-resolution networks with local attribution maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9199–9208,  
455 2021.  
456
- 457 [20] Jinjin Gu, Xianzheng Ma, Xiangtao Kong, Yu Qiao, and Chao Dong. Networks are slacking off: Under-  
458 standing generalization problem in image deraining. *Advances in Neural Information Processing Systems*,  
459 36, 2023.
- 460 [21] Rishi Hazra, Pedro Zuidberg Dos Martires, and Luc De Raedt. Saycanpay: Heuristic planning with large  
461 language models using learnable domain knowledge. In *Proceedings of the AAAI Conference on Artificial  
462 Intelligence*, volume 38, pages 20123–20133, 2024.
- 463 [22] Jinfan Hu, Jinjin Gu, Shiyao Yu, Fanghua Yu, Zheyuan Li, Zhiyuan You, Chaochao Lu, and Chao Dong.  
464 Interpreting low-level vision models with causal effect maps. *arXiv preprint arXiv:2407.19789*, 2024.
- 465 [23] Yuanming Hu, Hao He, Chenxi Xu, Baoyuan Wang, and Stephen Lin. Exposure: A white-box photo  
466 post-processing framework. *ACM Transactions on Graphics (TOG)*, 37(2):1–17, 2018.
- 467 [24] Jemin Hwangbo, Joonho Lee, Alexey Dosovitskiy, Dario Bellicoso, Vassilios Tsounis, Vladlen Koltun, and  
468 Marco Hutter. Learning agile and dynamic motor skills for legged robots. *Science Robotics*, 4(26):eaau5872,  
469 2019.
- 470 [25] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego  
471 de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b.  
472 *arXiv preprint arXiv:2310.06825*, 2023.
- 473 [26] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford,  
474 Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of  
475 experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 476 [27] Gu Jinjin, Cai Haoming, Chen Haoyu, Ye Xiaoxing, Jimmy S Ren, and Dong Chao. Pipal: a large-scale  
477 image quality assessment dataset for perceptual image restoration. In *Computer Vision—ECCV 2020:  
478 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pages 633–651.  
479 Springer, 2020.
- 480 [28] Daniel Kahneman. Thinking, fast and slow. *Farrar, Straus and Giroux*, 2011.
- 481 [29] Manzoor Ahmed Khan, Hesham El Sayed, Sumbal Malik, Talha Zia, Jalal Khan, Najla Alkaabi, and Henry  
482 Ignatious. Level-5 autonomous driving—are we there yet? a review of research literature. *ACM Computing  
483 Surveys (CSUR)*, 55(2):1–38, 2022.
- 484 [30] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language  
485 models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213,  
486 2022.
- 487 [31] Xiangtao Kong, Jinjin Gu, Yihao Liu, Wenlong Zhang, Xiangyu Chen, Yu Qiao, and Chao Dong. A  
488 preliminary exploration towards general image restoration. *arXiv preprint arXiv:2408.15143*, 2024.
- 489 [32] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- 490 [33] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich  
491 Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-  
492 intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- 493 [34] Chunyuan Li, Zhe Gan, Zhengyuan Yang, Jianwei Yang, Linjie Li, Lijuan Wang, Jianfeng Gao, et al.  
494 Multimodal foundation models: From specialists to general-purpose assistants. *Foundations and Trends®  
495 in Computer Graphics and Vision*, 16(1-2):1–214, 2024.
- 496 [35] Bill Yuchen Lin, Yicheng Fu, Karina Yang, Faeze Brahma, Shiyu Huang, Chandra Bhagavatula, Prithviraj  
497 Ammanabrolu, Yejin Choi, and Xiang Ren. Swiftsage: A generative agent with fast and slow thinking for  
498 complex interactive tasks. *Advances in Neural Information Processing Systems*, 36, 2024.
- 499 [36] Yunlong Lin, Zixu Lin, Haoyu Chen, Panwang Pan, Chenxin Li, Sixiang Chen, Yeying Jin, Wenbo Li, and  
500 Xinghao Ding. Jarvis: Elevating autonomous driving perception with intelligent image restoration. *arXiv  
501 preprint arXiv:2504.04158*, 2025.
- 502 [37] Yihao Liu, Anran Liu, Jinjin Gu, Zhipeng Zhang, Wenhao Wu, Yu Qiao, and Chao Dong. Discovering  
503 distinctive " semantics" in super-resolution networks. *arXiv preprint arXiv:2108.00406*, 2021.

- 504 [38] Yihao Liu, Hengyuan Zhao, Jinjin Gu, Yu Qiao, and Chao Dong. Evaluating the generalization ability of  
505 super-resolution networks. *IEEE Transactions on pattern analysis and machine intelligence*, 2023.
- 506 [39] Salma Abdel Magid, Zudi Lin, Donglai Wei, Yulun Zhang, Jinjin Gu, and Hanspeter Pfister. Texture-based  
507 error analysis for image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
508 and Pattern Recognition*, pages 2118–2127, 2022.
- 509 [40] Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam  
510 Singh, Xiangru Tang, Leandro Von Werra, and Shayne Longpre. Octopack: Instruction tuning code large  
511 language models. *arXiv preprint arXiv:2308.07124*, 2023.
- 512 [41] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang,  
513 Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with  
514 human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- 515 [42] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S  
516 Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual  
517 acm symposium on user interface software and technology*, pages 1–22, 2023.
- 518 [43] Tal Ridnik, Dedy Kredo, and Itamar Friedman. Code generation with alphacodium: From prompt  
519 engineering to flow engineering. *arXiv preprint arXiv:2401.08500*, 2024.
- 520 [44] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt:  
521 Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing  
522 Systems*, 36, 2024.
- 523 [45] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:  
524 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*,  
525 36, 2024.
- 526 [46] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,  
527 Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. A general reinforcement learning  
528 algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- 529 [47] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix,  
530 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation  
531 language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 532 [48] Haoran Wang, Yue Zhang, and Xiaosheng Yu. An overview of image caption generation methods.  
533 *Computational intelligence and neuroscience*, 2020(1):3062706, 2020.
- 534 [49] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,  
535 Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of  
536 Computer Science*, 18(6):186345, 2024.
- 537 [50] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-  
538 resolution with pure synthetic data. In *Proceedings of the IEEE/CVF international conference on computer  
539 vision*, pages 1905–1914, 2021.
- 540 [51] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
541 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural  
542 information processing systems*, 35:24824–24837, 2022.
- 543 [52] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt:  
544 Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- 545 [53] Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li, Wenxiu  
546 Sun, Qiong Yan, Guangtao Zhai, et al. Q-bench: A benchmark for general-purpose foundation models on  
547 low-level vision. *arXiv preprint arXiv:2309.14181*, 2023.
- 548 [54] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan  
549 Wang, Erli Zhang, Wenxiu Sun, et al. Q-align: Teaching llms for visual scoring via discrete text-defined  
550 levels. *arXiv preprint arXiv:2312.17090*, 2023.
- 551 [55] Haoning Wu, Hanwei Zhu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Annan  
552 Wang, Wenxiu Sun, Qiong Yan, et al. Towards open-ended visual quality comparison. *arXiv preprint  
553 arXiv:2402.16641*, 2024.

- 554 [56] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang,  
555 Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. *arXiv*  
556 *preprint arXiv:2309.07864*, 2023.
- 557 [57] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.  
558 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information*  
559 *Processing Systems*, 36, 2024.
- 560 [58] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React:  
561 Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- 562 [59] Zhiyuan You, Xin Cai, Jinjin Gu, Tianfan Xue, and Chao Dong. Teaching large language models to regress  
563 accurate image quality scores using score distribution. *arXiv preprint arXiv:2501.11561*, 2025.
- 564 [60] Zhiyuan You, Jinjin Gu, Zheyuan Li, Xin Cai, Kaiwen Zhu, Tianfan Xue, and Chao Dong. Descriptive  
565 image quality assessment in the wild. *arXiv preprint arXiv:2405.18842*, 2024.
- 566 [61] Zhiyuan You, Zheyuan Li, Jinjin Gu, Zhenfei Yin, Tianfan Xue, and Chao Dong. Depicting beyond  
567 scores: Advancing image quality assessment through multi-modal language models. *arXiv preprint*  
568 *arXiv:2312.08962*, 2023.
- 569 [62] Fanghua Yu, Jinjin Gu, Zheyuan Li, Jinfan Hu, Xiangtao Kong, Xintao Wang, Jingwen He, Yu Qiao, and  
570 Chao Dong. Scaling up to excellence: Practicing model scaling for photo-realistic image restoration in the  
571 wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
572 25669–25680, 2024.
- 573 [63] Ke Yu, Chao Dong, Liang Lin, and Chen Change Loy. Crafting a toolchain for image restoration by deep  
574 reinforcement learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
575 pages 2443–2452, 2018.
- 576 [64] MH Fazel Zarandi, Marzie Zarinbal, and Mina Izadi. Systematic image processing for diagnosing brain  
577 tumors: A type-ii fuzzy expert system approach. *Applied soft computing*, 11(1):285–294, 2011.
- 578 [65] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser:  
579 Residual learning of deep cnn for image denoising. *IEEE transactions on image processing*, 26(7):3142–  
580 3155, 2017.
- 581 [66] Ruofan Zhang, Jinjin Gu, Haoyu Chen, Chao Dong, Yulun Zhang, and Wenming Yang. Crafting training  
582 degradation distribution for the accuracy-generalization trade-off in real-world super-resolution. In  
583 *International conference on machine learning*, pages 41078–41091. PMLR, 2023.
- 584 [67] Wenlong Zhang, Xiaohui Li, Xiangyu Chen, Yu Qiao, Xiao-Ming Wu, and Chao Dong. Seal: A framework  
585 for systematic evaluation of real-world super-resolution. *arXiv preprint arXiv:2309.03020*, 2023.
- 586 [68] Wenlong Zhang, Guangyuan Shi, Yihao Liu, Chao Dong, and Xiao-Ming Wu. A closer look at blind  
587 super-resolution: Degradation models, baselines, and performance upper bounds. In *Proceedings of the*  
588 *IEEE/CVF conference on computer vision and pattern recognition*, pages 527–536, 2022.
- 589 [69] Kaiwen Zhu, Jinjin Gu, Zhiyuan You, Yu Qiao, and Chao Dong. An intelligent agentic system for complex  
590 image restoration problems. *arXiv preprint arXiv:2410.17809*, 2024.
- 591 [70] Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei  
592 Lu, Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world environments  
593 via large language models with text-based knowledge and memory. *arXiv preprint arXiv:2305.17144*,  
594 2023.

## 595 Appendix

### 596 A More Problems Demanding Further Study

#### 597 A.1 Image Quality Assessment and Content Analysis

598 Regardless of how we design our cognitive architecture, certain fundamental capabilities are in-  
599 dispensable. Among these, the recognition, analysis, and evaluation of image content and quality  
600 are essential. In Figure 4, the Perception and Reflection stages are related to this capability; they

601 serve as the “eyes” of the image processing agentic system. Historically, image content recognition  
602 [48] and image quality assessment [27, 18] have been independent research fields, separate from  
603 image processing models, each with its own methods and objectives. However, from the research  
604 perspective of agentic systems, we impose higher demands on them.

605 For image content recognition, we need to consider its robustness under different image qualities  
606 and special circumstances. The recognized content must be designed with finer granularity to meet  
607 the specific needs of image processing, rather than being overly abstract like high-level vision tasks.  
608 Regarding image quality assessment, we cannot limit ourselves to evaluating a single “score.” Models  
609 need to be more intelligent, performing fine-grained image quality analysis, determining types of  
610 distortion, and assessing the quality of intermediate results. The decisions of the entire agentic system  
611 largely depend on the accuracy and intelligence of this pair of “eyes.”

612 Thanks to the development of multi-modal language models, a number of tools have now begun  
613 to demonstrate this capability [61, 60, 53, 54, 55]. Utilizing language models, these methods can  
614 describe image content, analyze image quality, evaluate the pros and cons of different image qualities,  
615 and provide judgments based on quality. These methods have already been applied in early research  
616 on agent-based image processing systems, showcasing their potential. However, the intelligence and  
617 accuracy of these methods still have a significant gap compared to large-scale practical applications.

## 618 **A.2 Knowledge Acquisition and Infusion**

619 After discussing the “eyes,” let’s turn to the “brain.” As previously analyzed, the planning and  
620 decision-making abilities in intelligent agentic systems mainly stem from language models, which  
621 base their decisions on general knowledge learned from large volumes of text training data. Only  
622 when a language model has encountered problems and knowledge related to image processing during  
623 training can it be expected to make accurate judgments in the system; otherwise, the language model  
624 may struggle to provide reliable predictions. However, pre-trained language models usually contain  
625 only the most basic knowledge. If we want an image processing system to perform more specific and  
626 precise tasks, we need to supply the language model with the necessary knowledge. This involves  
627 two issues: the acquisition of knowledge and the injection of knowledge.

628 Firstly, acquiring such knowledge is non-trivial, and the method of injecting this knowledge into  
629 the system depends heavily on how it is represented. Image processing involves not only a large  
630 amount of conceptual and systematic knowledge but also a wealth of experience-based and case-based  
631 knowledge. This kind of knowledge is difficult to abstract into rules and usually exists in the form of  
632 case studies. Early attempts mainly employed two methods. Chen et al. [9] collected a series of input  
633 images along with corresponding instruction-output training data to implicitly carry a large amount  
634 of knowledge and information. The trained model then has the ability to handle similar problems.  
635 However, the drawbacks of this method are evident: firstly, collecting a large amount of high-quality  
636 training data requires substantial resources and is both costly and difficult. Secondly, the model  
637 lacks scalability; adding new knowledge requires retraining the model. Additionally, fine-tuning the  
638 language model may compromise its general capabilities.

639 In contrast, Zhu et al. [69] rely on the reasoning ability of an unmodified language model and provide  
640 a reference “manual.” This method of supplying knowledge is known as Retrieval-Augmented  
641 Generation [33]. They first use the language model to summarize a large amount of scattered case  
642 information into knowledge that can be described linguistically. When solving actual problems,  
643 they provide the relevant content together. The language model utilizes its zero-shot learning and  
644 contextual inference capabilities to complete tasks based on the provided information. The advantage  
645 of this method is that it does not fine-tune the model, avoiding the loss of general performance, and  
646 the model still possesses strong reasoning and understanding abilities. However, its drawbacks lie  
647 in the difficulty of accurately describing professional knowledge in language. Moreover, quickly  
648 retrieving relevant information from a vast amount of knowledge is not easy, and it also places high  
649 demands on the design of the cognitive architecture.

650 In this area, we currently have only very preliminary results. The exploration, acquisition, representa-  
651 tion, and injection of prior knowledge in image processing will become the core research topics of  
652 image processing agent systems.

### 653 **A.3 Human-Computer Interaction in Agentic System**

654 Agentic systems' multi-step operational paradigms, randomness, and natural language interfaces  
655 introduce new challenges in human-computer interaction. Currently, the primary way to interact  
656 with intelligent agent systems is through "chatting," where the system communicates its thoughts  
657 and actions in a conversational manner. However, we need new interaction methods to meet higher  
658 demands.

659 We must provide users with visibility into what the agent is doing by displaying all the steps it takes,  
660 allowing users to observe and understand the ongoing processes. Simultaneously, users should be  
661 able to give the agent more fine-grained and explicit instructions to control its behavior more precisely.  
662 Moreover, users should not only see what is happening but also have the ability to correct the agent.  
663 If they discover that the agent made an incorrect choice at step four (out of ten), they should be able  
664 to return to that step, correct the agent in some manner, and then proceed with the execution. The  
665 ultimate goal is to achieve collaboration between the agent and the user, enabling them to complete  
666 tasks together effectively.

### 667 **A.4 Exploitative Learning, Self-Evolution & Creativity**

668 The development of agentic systems has opened up new horizons in the fields of exploitative learning,  
669 self-evolution, and creativity. These concepts are crucial for advancing intelligent systems, enabling  
670 them to autonomously adapt, improve, and innovate over time without explicit human intervention,  
671 achieving higher levels of automation and agency as depicted in Figure 3.

672 Exploitative learning refers to the agent itself taking the initiative to determine the methods and  
673 content of knowledge acquisition within certain limits. The work of Chen et al. [10] embodies the  
674 prototype of this idea: they presented many experimental results to the agent, and the agent selected  
675 valuable content from them to learn. In some cases, the agent could even take some unconventional  
676 actions to acquire new knowledge through interaction with the world.

677 Self-evolution is the agent's ability to develop its own algorithms and strategies over time. This not  
678 only involves learning from data but also enables the agent to continuously improve itself based on  
679 its processing results, learning from past cases. Through iterative self-assessment and refinement, the  
680 agent gradually enhances its performance and may even modify its underlying processes to better  
681 adapt to changing environments or objectives.

682 Creativity in agent systems goes beyond mere problem-solving; it includes generating new ideas,  
683 methods, or outputs that are both original and valuable. This involves not only developing unique  
684 approaches to tackle complex image processing challenges that standard algorithms cannot handle,  
685 but also creatively generating content such as artistic transformations, stylizations, or entirely new  
686 visual effects.

687 These are grand visions under higher levels of agency and autonomy. At this stage, exploration of  
688 these issues is still quite limited. This paper serves only as an introduction to envisioning these higher  
689 levels of intelligence.

## 690 **NeurIPS Paper Checklist**

691 The checklist is designed to encourage best practices for responsible machine learning research,  
692 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
693 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
694 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
695 towards the page limit.

696 Please read the checklist guidelines carefully for information on how to answer these questions. For  
697 each question in the checklist:

- 698 • You should answer [Yes], [No], or [NA].
- 699 • [NA] means either that the question is Not Applicable for that particular paper or the  
700 relevant information is Not Available.
- 701 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

702 **The checklist answers are an integral part of your paper submission.** They are visible to the  
703 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
704 (after eventual revisions) with the final version of your paper, and its final version will be published  
705 with the paper.

706 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
707 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
708 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
709 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
710 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
711 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
712 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
713 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
714 please point to the section(s) where related material for the question can be found.

715 **IMPORTANT, please:**

- 716 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 717 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 718 • **Do not modify the questions and only use the provided macros for your answers.**

### 719 **1. Claims**

720 Question: Do the main claims made in the abstract and introduction accurately reflect the  
721 paper’s contributions and scope?

722 Answer: [Yes]

723 Justification: The main claims made in the abstract and introduction accurately reflect the  
724 paper’s contributions and scope.

725 Guidelines:

- 726 • The answer NA means that the abstract and introduction do not include the claims  
727 made in the paper.
- 728 • The abstract and/or introduction should clearly state the claims made, including the  
729 contributions made in the paper and important assumptions and limitations. A No or  
730 NA answer to this question will not be perceived well by the reviewers.
- 731 • The claims made should match theoretical and experimental results, and reflect how  
732 much the results can be expected to generalize to other settings.
- 733 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
734 are not attained by the paper.

### 735 **2. Limitations**

736 Question: Does the paper discuss the limitations of the work performed by the authors?

737 Answer: [Yes]

738 Justification: The paper discusses the limitations of the work.

739 Guidelines:

- 740 • The answer NA means that the paper has no limitation while the answer No means that
- 741 the paper has limitations, but those are not discussed in the paper.
- 742 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 743 • The paper should point out any strong assumptions and how robust the results are to
- 744 violations of these assumptions (e.g., independence assumptions, noiseless settings,
- 745 model well-specification, asymptotic approximations only holding locally). The authors
- 746 should reflect on how these assumptions might be violated in practice and what the
- 747 implications would be.
- 748 • The authors should reflect on the scope of the claims made, e.g., if the approach was
- 749 only tested on a few datasets or with a few runs. In general, empirical results often
- 750 depend on implicit assumptions, which should be articulated.
- 751 • The authors should reflect on the factors that influence the performance of the approach.
- 752 For example, a facial recognition algorithm may perform poorly when image resolution
- 753 is low or images are taken in low lighting. Or a speech-to-text system might not be
- 754 used reliably to provide closed captions for online lectures because it fails to handle
- 755 technical jargon.
- 756 • The authors should discuss the computational efficiency of the proposed algorithms
- 757 and how they scale with dataset size.
- 758 • If applicable, the authors should discuss possible limitations of their approach to
- 759 address problems of privacy and fairness.
- 760 • While the authors might fear that complete honesty about limitations might be used by
- 761 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
- 762 limitations that aren't acknowledged in the paper. The authors should use their best
- 763 judgment and recognize that individual actions in favor of transparency play an impor-
- 764 tant role in developing norms that preserve the integrity of the community. Reviewers
- 765 will be specifically instructed to not penalize honesty concerning limitations.

### 766 3. Theory assumptions and proofs

767 Question: For each theoretical result, does the paper provide the full set of assumptions and

768 a complete (and correct) proof?

769 Answer: [NA]

770 Justification: No theory.

771 Guidelines:

- 772 • The answer NA means that the paper does not include theoretical results.
- 773 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 774 referenced.
- 775 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 776 • The proofs can either appear in the main paper or the supplemental material, but if
- 777 they appear in the supplemental material, the authors are encouraged to provide a short
- 778 proof sketch to provide intuition.
- 779 • Inversely, any informal proof provided in the core of the paper should be complemented
- 780 by formal proofs provided in appendix or supplemental material.
- 781 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 782 4. Experimental result reproducibility

783 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

784 perimental results of the paper to the extent that it affects the main claims and/or conclusions

785 of the paper (regardless of whether the code and data are provided or not)?

786 Answer: [NA]

787 Justification: No experiments

788 Guidelines:

- 789 • The answer NA means that the paper does not include experiments.

- 790 • If the paper includes experiments, a No answer to this question will not be perceived  
791 well by the reviewers: Making the paper reproducible is important, regardless of  
792 whether the code and data are provided or not.
- 793 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
794 to make their results reproducible or verifiable.
- 795 • Depending on the contribution, reproducibility can be accomplished in various ways.  
796 For example, if the contribution is a novel architecture, describing the architecture fully  
797 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
798 be necessary to either make it possible for others to replicate the model with the same  
799 dataset, or provide access to the model. In general, releasing code and data is often  
800 one good way to accomplish this, but reproducibility can also be provided via detailed  
801 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
802 of a large language model), releasing of a model checkpoint, or other means that are  
803 appropriate to the research performed.
- 804 • While NeurIPS does not require releasing code, the conference does require all submis-  
805 sions to provide some reasonable avenue for reproducibility, which may depend on the  
806 nature of the contribution. For example
  - 807 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
808 to reproduce that algorithm.
  - 809 (b) If the contribution is primarily a new model architecture, the paper should describe  
810 the architecture clearly and fully.
  - 811 (c) If the contribution is a new model (e.g., a large language model), then there should  
812 either be a way to access this model for reproducing the results or a way to reproduce  
813 the model (e.g., with an open-source dataset or instructions for how to construct  
814 the dataset).
  - 815 (d) We recognize that reproducibility may be tricky in some cases, in which case  
816 authors are welcome to describe the particular way they provide for reproducibility.  
817 In the case of closed-source models, it may be that access to the model is limited in  
818 some way (e.g., to registered users), but it should be possible for other researchers  
819 to have some path to reproducing or verifying the results.

## 820 5. Open access to data and code

821 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
822 tions to faithfully reproduce the main experimental results, as described in supplemental  
823 material?

824 Answer: [NA]

825 Justification: No experiments.

826 Guidelines:

- 827 • The answer NA means that paper does not include experiments requiring code.
- 828 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
829 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 830 • While we encourage the release of code and data, we understand that this might not be  
831 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
832 including code, unless this is central to the contribution (e.g., for a new open-source  
833 benchmark).
- 834 • The instructions should contain the exact command and environment needed to run to  
835 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
836 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 837 • The authors should provide instructions on data access and preparation, including how  
838 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 839 • The authors should provide scripts to reproduce all experimental results for the new  
840 proposed method and baselines. If only a subset of experiments are reproducible, they  
841 should state which ones are omitted from the script and why.
- 842 • At submission time, to preserve anonymity, the authors should release anonymized  
843 versions (if applicable).

- 844 • Providing as much information as possible in supplemental material (appended to the  
845 paper) is recommended, but including URLs to data and code is permitted.

846 **6. Experimental setting/details**

847 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
848 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
849 results?

850 Answer: [NA]

851 Justification: No experiments

852 Guidelines:

- 853 • The answer NA means that the paper does not include experiments.  
854 • The experimental setting should be presented in the core of the paper to a level of detail  
855 that is necessary to appreciate the results and make sense of them.  
856 • The full details can be provided either with the code, in appendix, or as supplemental  
857 material.

858 **7. Experiment statistical significance**

859 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
860 information about the statistical significance of the experiments?

861 Answer: [NA]

862 Justification: No experiments

863 Guidelines:

- 864 • The answer NA means that the paper does not include experiments.  
865 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
866 dence intervals, or statistical significance tests, at least for the experiments that support  
867 the main claims of the paper.  
868 • The factors of variability that the error bars are capturing should be clearly stated (for  
869 example, train/test split, initialization, random drawing of some parameter, or overall  
870 run with given experimental conditions).  
871 • The method for calculating the error bars should be explained (closed form formula,  
872 call to a library function, bootstrap, etc.)  
873 • The assumptions made should be given (e.g., Normally distributed errors).  
874 • It should be clear whether the error bar is the standard deviation or the standard error  
875 of the mean.  
876 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
877 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
878 of Normality of errors is not verified.  
879 • For asymmetric distributions, the authors should be careful not to show in tables or  
880 figures symmetric error bars that would yield results that are out of range (e.g. negative  
881 error rates).  
882 • If error bars are reported in tables or plots, The authors should explain in the text how  
883 they were calculated and reference the corresponding figures or tables in the text.

884 **8. Experiments compute resources**

885 Question: For each experiment, does the paper provide sufficient information on the com-  
886 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
887 the experiments?

888 Answer: [NA]

889 Justification: No experiments.

890 Guidelines:

- 891 • The answer NA means that the paper does not include experiments.  
892 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
893 or cloud provider, including relevant memory and storage.

- 894
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 895
- 896
- 897
- 898

## 899 9. Code of ethics

900 Question: Does the research conducted in the paper conform, in every respect, with the  
901 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

902 Answer: [Yes]

903 Justification: tThe research conducted in the paper conform, in every respect, with the  
904 NeurIPS Code of Ethics.

905 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 906
- 907
- 908
- 909
- 910

## 911 10. Broader impacts

912 Question: Does the paper discuss both potential positive societal impacts and negative  
913 societal impacts of the work performed?

914 Answer: [NA]

915 Justification: tThere is no societal impact of the work performed.

916 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
  - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 917
- 918
- 919
- 920
- 921
- 922
- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936
- 937
- 938

## 939 11. Safeguards

940 Question: Does the paper describe safeguards that have been put in place for responsible  
941 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
942 image generators, or scraped datasets)?

943 Answer: [NA]

944 Justification: tThe paper poses no such risks.

945 Guidelines:

- 946 • The answer NA means that the paper poses no such risks.
- 947 • Released models that have a high risk for misuse or dual-use should be released with
- 948 necessary safeguards to allow for controlled use of the model, for example by requiring
- 949 that users adhere to usage guidelines or restrictions to access the model or implementing
- 950 safety filters.
- 951 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 952 should describe how they avoided releasing unsafe images.
- 953 • We recognize that providing effective safeguards is challenging, and many papers do
- 954 not require this, but we encourage authors to take this into account and make a best
- 955 faith effort.

## 956 12. Licenses for existing assets

957 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
958 the paper, properly credited and are the license and terms of use explicitly mentioned and  
959 properly respected?

960 Answer: [NA]

961 Justification: The paper does not use existing assets.

962 Guidelines:

- 963 • The answer NA means that the paper does not use existing assets.
- 964 • The authors should cite the original paper that produced the code package or dataset.
- 965 • The authors should state which version of the asset is used and, if possible, include a
- 966 URL.
- 967 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 968 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 969 service of that source should be provided.
- 970 • If assets are released, the license, copyright information, and terms of use in the
- 971 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)
- 972 has curated licenses for some datasets. Their licensing guide can help determine the
- 973 license of a dataset.
- 974 • For existing datasets that are re-packaged, both the original license and the license of
- 975 the derived asset (if it has changed) should be provided.
- 976 • If this information is not available online, the authors are encouraged to reach out to
- 977 the asset's creators.

## 978 13. New assets

979 Question: Are new assets introduced in the paper well documented and is the documentation  
980 provided alongside the assets?

981 Answer: [NA]

982 Justification: The paper does not release new assets.

983 Guidelines:

- 984 • The answer NA means that the paper does not release new assets.
- 985 • Researchers should communicate the details of the dataset/code/model as part of their
- 986 submissions via structured templates. This includes details about training, license,
- 987 limitations, etc.
- 988 • The paper should discuss whether and how consent was obtained from people whose
- 989 asset is used.
- 990 • At submission time, remember to anonymize your assets (if applicable). You can either
- 991 create an anonymized URL or include an anonymized zip file.

## 992 14. Crowdsourcing and research with human subjects

993 Question: For crowdsourcing experiments and research with human subjects, does the paper  
994 include the full text of instructions given to participants and screenshots, if applicable, as  
995 well as details about compensation (if any)?

996 Answer: [NA]

997 Justification: The paper does not involve crowdsourcing nor research with human subjects.

998 Guidelines:

- 999 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1000 human subjects.
- 1001 • Including this information in the supplemental material is fine, but if the main contribu-
- 1002 tion of the paper involves human subjects, then as much detail as possible should be
- 1003 included in the main paper.
- 1004 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 1005 or other labor should be paid at least the minimum wage in the country of the data
- 1006 collector.

1007 **15. Institutional review board (IRB) approvals or equivalent for research with human**

1008 **subjects**

1009 Question: Does the paper describe potential risks incurred by study participants, whether

1010 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

1011 approvals (or an equivalent approval/review based on the requirements of your country or

1012 institution) were obtained?

1013 Answer: [NA]

1014 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1015 Guidelines:

- 1016 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1017 human subjects.
- 1018 • Depending on the country in which research is conducted, IRB approval (or equivalent)
- 1019 may be required for any human subjects research. If you obtained IRB approval, you
- 1020 should clearly state this in the paper.
- 1021 • We recognize that the procedures for this may vary significantly between institutions
- 1022 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 1023 guidelines for their institution.
- 1024 • For initial submissions, do not include any information that would break anonymity (if
- 1025 applicable), such as the institution conducting the review.

1026 **16. Declaration of LLM usage**

1027 Question: Does the paper describe the usage of LLMs if it is an important, original, or

1028 non-standard component of the core methods in this research? Note that if the LLM is used

1029 only for writing, editing, or formatting purposes and does not impact the core methodology,

1030 scientific rigorousness, or originality of the research, declaration is not required.

1031 Answer: [Yes]

1032 Justification: We use LLM to proofreading.

1033 Guidelines:

- 1034 • The answer NA means that the core method development in this research does not
- 1035 involve LLMs as any important, original, or non-standard components.
- 1036 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
- 1037 for what should or should not be described.