PROACTIVE AGENTS FOR MULTI-TURN TEXT-TO-IMAGE GENERATION UNDER UNCERTAINTY

Anonymous authors

Paper under double-blind review

ABSTRACT

User prompts for generative AI models are often underspecified or open-ended, which may lead to sub-optimal responses. This prompt underspecification problem is particularly evident in text-to-image (T2I) generation, where users commonly struggle to articulate their precise intent. This disconnect between the user's vision and the model's interpretation often forces users to painstakingly and repeatedly refine their prompts. To address this, we propose a design for proactive T2I agents equipped with an interface to actively ask clarification questions when uncertain, and present their understanding of user intent as an interpretable *belief graph* that a user can edit. We build simple prototypes for such agents and verify their effectiveness through both human studies and automated evaluation. We observed that at least 90% of human subjects found these agents and their belief graphs helpful for their T2I workflow. Moreover, we use a scalable automated evaluation approach using two agents, one with a ground truth image and the other tries to ask as few questions as possible to align with the ground truth. On DesignBench, a benchmark we created for artists and designers, the COCO dataset (Lin et al., 2014) and ImageInWords (Garg et al., 2024), we observed that these T2I agents were able to ask informative questions and elicit crucial information to achieve successful alignment with at least 2 times higher VQAScore (Lin et al., 2024) than the standard single-turn T2I generation. Demo: https://youtu.be/HPqJ4xPRnto





054 1 INTRODUCTION

A fundamental challenge in the development of AI agents is how to foster effective and efficient multiturn communication and collaboration with human users to achieve user-defined goals, especially when faced with the common issue of vague or incomplete instructions from humans. We focus specifically on text-to-image (T2I) generation, where recent advancements (Baldridge et al., 2024; Betker et al., 2023; Podell et al., 2023; Yu et al., 2023) have enabled the creation of stunning images from complex text descriptions. However, users often struggle to describe the image they would like to generate in a way that T2I systems can fully understand. This leads to unsatisfactory results and repeated iterations of prompts.

064 The prompt underspecification problem arises from the inherent ambiguity of natural language, the different assumptions that humans make and the vast space of potential images that can be generated 065 from a single prompt (Hutchinson et al., 2022). Imagine a prompt generate an image of a rabbit next 066 to a cat. This seemingly simple prompt leaves many important aspects underspecified: What kind of 067 rabbit? What color is the cat? What is their relative positions? What is the background? While a 068 T2I model can generate an image with a rabbit and a cat in it, it is unlikely that the image captures 069 the specific details a specific user has in mind. For example, people in Holland might assume it is common for rabbits to have lop ears, but people in New England might expect to see cottontail rabbit 071 with straight ears. The combination of all these factors can lead to a frustrating cycle of trial-and-error, 072 with the user repeatedly refining their prompt in an attempt to steer the model towards the desired 073 output (Vodrahalli & Zou, 2024; Huang et al., 2024; Sun & Guo, 2023).

Instead of relying on passive T2I models that simply generate images based on potentially vague user instructions, we pursue a quest for agency in T2I generation. The T2I agents should actively engage with human users to provide a collaborative and interactive experience for image creation. We envision that these T2I agents will be able to (1) express and visualize their beliefs and uncertainty about user intents, (2) allow human users to directly control their beliefs beyond just text descriptions, and (3) proactively seek clarification from the human user to iteratively align their understanding with what the human user intends to generate.

081 In this work, we develop simple prototypes of such agents. At the core of those agent prototypes, we 082 build in a graph-based symbolic belief state, named *belief graph*, for agents to understand its own 083 uncertainty about possible entities (e.g., rabbit) that might appear in the image, attributes of entities 084 (e.g., rabbit's color), relations between entities and so on. Given a user prompt, we use an LLM 085 and constrain its generation to the graph structure of beliefs, which include probability estimates on the appearance of entities and the possible values for attributes and relations. Figure 1 illustrates 087 the interface and features of the prototypes. In particular, the agent can ask questions based on its uncertainty. For example, a very simple strategy is to find the most uncertain attribute of an entity 880 (e.g., rabbit's color) and use an LLM to phrase a question about the attribute (e.g., What is the color 089 of the rabbit?). The agent can also guide users to directly edit items in the graph. 090

To evaluate the utility of our agent prototypes, we conduct both human studies and automatic evaluations. The human studies aim to understand how helpful simple T2I agents can be, and evaluate how good the agents' questions are. We develop automatic evaluation pipelines to assess the effectiveness and efficiency of the T2I agents when interacting with simulated users with underspecified prompts answering questions based on their pre-fixed intents.

096 We found that over 90% human subjects expect proactive clarifications to be helpful, and 58% think this question asking feature of agents could deliver value to their work very soon, or immediately. We 098 create a hand-curated benchmark called DesignBench which contains aesthetic scenes with multiple entities and interactions between entities; it also contains both a short and long caption. DesignBench also features diversity between photo-realism, animation and multiple styles allowing a robust testing 100 with the use case of artists and designers in mind. This benchmark will be released with this paper. 101 We run automatic evaluations on both the COCO dataset (Lin et al., 2014) and DesignBench. We 102 found that our agents can achieve at least 2 times higher VQAScore (Lin et al., 2024) than the 103 traditional single-turn T2I generation within 5 turns of interaction. 104

Our contributions: (1) the first interpretable and controllable belief graph used for T2I, (2) novel
 design and prototypes for T2I agents that adaptively ask clarification questions and present belief
 graphs; (3) a new automatic evaluation pipeline with simulated users to assess question-asking skills
 of T2I agents; and (4) DesignBench: a new T2I agent benchmark. Appendix A details the novelty.

108 2 RELATED WORK

109

From the very outset of artificial intelligence, a core challenge has been to develop intelligent agents capable of representing knowledge and taking actions to acquire knowledge necessary for achieving their goals (McCarthy & Hayes, 1969; Minsky, 1974; Moore, 1985; Nilsson, 2009; Russell & Norvig, 2016). Our work is an attempt to address this challenge for intelligent T2I agents.

In machine learning and statistics, efficient data acquisition has been extensively studied for many problems, including active learning (Settles, 2009), Bayesian optimization (Garnett, 2023), reinforcement learning (Kaelbling et al., 1996; Sutton, 2018) and experimental design (Chaloner & Verdinelli, 1995; Kirk, 2009). We reckon that T2I agents should also be capable of actively seeking important information from human users to quickly reduce uncertainty (Wang et al., 2024b) and generate satisfying images. In §D, we detail the implementation of action selection strategies for our T2I agents.

In human-computer interaction, researchers have been extensively studying how to best enable
Human-AI interaction especially from user experience perspectives (Norman, 1994; Höök, 2000;
Amershi et al., 2019; Cai et al., 2019; Viégas & Wattenberg, 2023; Chen et al., 2024; Yang et al.,
2020; Kim et al., 2023). Interface design for AI is becoming increasingly challenging due to the
lack of transparency (Viégas & Wattenberg, 2023; Chen et al., 2024), uncertainty about AI capability
and complex outputs (Yang et al., 2020). We aim to build user-friendly agents, and an indispensable
component is their interface to enable them to effectively act and observe, as detailed in §G.

Interpretebaility. Surfacing an agent's belief overlaps with interpretability as both aim to understand model or agent's internal. Some methods leverage LLM's natural language interface to surface their reasoning (e.g., chain of thought (Wei et al., 2023a)), sometime interactively (Wang et al., 2024a).
While these approaches makes accessible explanations, whether the explanations represents truth has been questioned (Lanham et al., 2023; Wei et al., 2023b; Chen et al., 2023). Some studies indicate explanations generated by the LLMs may not entail the models' predictions nor be factually grounded in the input, even on simple tasks with extractive explanations (Ye & Durrett, 2022).

Text-to-Image (T2I) generation. Text-to-image prompts can be ambiguous, subjective (Hutchinson et al., 2022), or challenging to represent visually (Wiles et al., 2024). Different users often have distinct requirements for image generation, including personalization (Wei et al., 2024), style constraints (Wang et al., 2023), and individual interpretations (Yin et al., 2019). To create images that better align with users' specific needs and interpretations, it is essential to actively communicate and interact with the user to understand the user's intent.

141 Multi-turn T2I. Current multi-turn T2I systems typically focus on multi-turn user instructions. 142 Huang et al. (2024); Sun & Guo (2023) propose multi-modal interactive dialogue systems which passively respond to user's natural language instructions. Mini DALL·E 3 (Lai et al., 2023) builds an 143 144 interactive T2I framework with an LLM in the loop to have a dialogue with the user via text chat and improve image generation and editing based on the entire conversation. Vodrahalli & Zou (2024) 145 collected and analyzed a dataset of human-AI interactions where users iteratively refine prompts for 146 T2I models to generate images similar to goal images (goal images are only visible to users). This 147 may require users to actively try prompts to understand model behaviors. On the contrary, our work 148 aims to reduce the burden on the user by actively asking questions to understand user intents. 149

A core challenge in multi-turn T2I is consistency (Cheng et al., 2024a;b; Zeqiang et al., 2023). Hu et al. (2024) introduce Instruct-Imagen, which is a model that follows complex multi-modal instructions. AudioStudio (Cheng et al., 2024a) is a multi-turn T2I framework aimed at subject consistencies while generating diverse and coherent images. These consistency improvement methods can be integrated into T2I agents but it is beyond the scope of this work. Our key focus is on the sequential decision making capability of agents to elicit user intents.

156 157

158

3 BACKGROUND

159 The belief graph in our work is closely related to symbolic world representations.

World states. In classical AI, researchers use symbolic representations to describe the world state (McCarthy & Hayes, 1969; Minsky, 1974; 1988; Pasula et al., 2007; Kaelbling & Lozano-Pérez, 2011). For example, in the blocks world (Ginsberg & Smith, 1988; Gupta & Nau, 1992; Alkhazraji

¹⁶² et al., 2020), a state can be

164

$is_block(a) \land is_red(a) \land on_table(a) \land is_block(b) \land is_blue(b) \land on(b,a),$

describing that there are a red block and a blue block, referred to as a and b, block a is on a table, and block b is on a. Such world states must include **entities** (e.g., a and b), their **attributes** (e.g., position on_table, characteristics is_block) and **relations** (e.g., on(b, a)) which are critical for enabling a robot to know and act in the world.

In linguistics, Davidson (1965; 1967b;a) introduce logic-based formalisms of meanings of sentences. The semantics of a sentence is decomposed to a set of atomic propositions, such that no propositions can be added or removed from the set to represent the meaning of the sentence. (Cho et al., 2023) propose Davidsonian Scene Graph (DSG) which represent an image description as a set of atomic propositions (and corresponding questions about each proposition) to evaluate T2I alignment.

We borrow the same concept as symbolic world representations and scene graphs, except that the agent needs to represent an imaginary world. The image generation problem can be viewed as taking a picture of the imaginary world. The world state should include all entities that are in the picture, together with their attributes and relations.

Belief states. Term "belief state" (Nilsson, 1986; Kaelbling et al., 1998) has been used to describe a distribution over states. E.g., for block *a*, we might have $p(on_table(a)) = 0.5$ and $p(\neg on_table(a)) = 0.5$, which means the agent is unsure whether the block is on a table. To represent the T2I agent's belief on which image to generate, we need to consider the distribution over all possible "worlds" in which the picture can be taken. This distribution can be described by the probabilities that an entity appears in the picture, an attribute gets assigned a certain value, etc.

184 185

187

4 PROACTIVE T2I AGENT DESIGN

We provide high-level principles and design that guide our agent how to behave and interact with users to generate desired images from text through multi-turn interactions. The goal of the agent is to generate images that match the user's intended image as closely as possible with minimal back-andforth, particularly in cases with underspecified prompts and the agent needs to gather information proactively. This requires a decision strategy on information gathering to trade off between the cost of interactions and the quality of generated images. The formal problem definition can be found in §B.

We equip the agent with the ability to gather information in two ways: ask clarification questions (§4.1) and express its uncertainty and understanding in a way that users can edit (§4.2). Once a piece of information is collected from a user, the agent also need to update its questions and uncertainty (§4.3). To enable all these agent behaviors, we need to situate the agent in an interface to effectively communicate with users (§G). In the following, we introduce the design of the above components under the interface, to ensure information efficiency for T2I generation.

200 4.1 What kind of questions should be asked?

We explain considerations in question asking and examples of strategies in this section.

203 4.1.1 PRINCIPLES

204 We identify the following principles for an agent to ask the user questions about the underspecified 205 prompt and their intended image: (i) Relevance: The question should be based on the user prompt. 206 (ii) Uncertainty Reduction: The question should aim to reduce the agent's uncertainty about the 207 attributes and contents of the image, the objects, the spatial layout, and the style. (iii) Easy-to-Answer: 208 The question should be as concise and direct as possible to ensure it is not too difficult for the user to 209 answer. (iv) No Redundancy: The question should not collect information present in the history of 210 interactions with the user. The Relevance and No Redundancy principles are self-explanatory, we 211 detail the other two principles below.

- The Uncertainty Reduction principle aims to let agent elicit information about various characteristics of the desired image, which the agent is unsure of.
- First, the agent needs to know what characteristics of images are important. Some examples include: (i) Attributes of the subjects, such as breed, size, or color, with questions like *What kind of rabbit*?

What color is the cat?; (ii) Spatial relationships between the subjects, such as proximity and relative position (*Are the rabbit and cat close to each other? Are they facing each other?*); (iii) Background information, such as location, style and time of day (*Are they in a park or at home?*); and (iv) Implicit entities that might not be explicitly mentioned in the initial prompt but are relevant to the user's vision (*Are there any other animals or people present?*).

Second, the agent needs to know its own uncertainty about those characteristics. In the agent's belief, the uncertainty is explicit. One strategy is to form questions about the image characteristics that the agent is most uncertain about. We discuss more in §D.2.

Third, the agent needs to update its own uncertainty once the user gives a response to its question (a.k.a. transition in §4.3). Then, it can construct questions again based on its updated uncertainty estimates. This iterative clarification process allows the agent to progressively refine its understanding of the user's intent and generate an image that more accurately reflects their desired output.

The Easy-to-Answer principle aims to reduce users' effort to respond to questions. One way is
 to have the agent provide some answer options, where options are what the agent believes likely to
 appear. E.g., *What color is the cat? (a) Black (b) Brown (c) Orange (d) Other (please specify).*

4.1.2 EXAMPLES OF QUESTION-ASKING STRATEGIES

Given the agent belief constructed from the user prompt (more details in §4.2), several basic approaches can be employed following the above principles. We construct simple agents with the following strategies, which are implemented and used in our experiments.

• Ag1 (§D.5): Rule-based question generation, which leverages predefined rules or heuristics to identify salient attributes, entities, or relationships that require clarification. For example, an LLM could be used to estimate the importance and likelihood of different components within the belief, and a heuristic could be applied to prioritize the most crucial elements for questioning.

Ag2 (§D.6): Belief-guided question generation, which involves using natural language to represent the current understanding encapsulated in the belief. This representation, along with the conversation history, is provided as input to an LLM, guiding it to generate clarification questions.

• Ag3 (§D.7): Direct question generation, which write the above question-asking principles in a prompt for an LLM to generate a question.

245 246 4.2 INTERACTING WITH THE USER BASED ON AGENT BELIEFS

The Uncertainty Reduction principle inspires the usage of belief graphs for the agent to directly
express uncertainty, in addition to reflecting uncertainty through questions. Instead of using hardcoded
symbols in classic belief representations (Fikes & Nilsson, 1971) described in §3, we employ LLMs
to generate names and values for entities, attributes and relations. As a result, this belief construction
method can generalize across any prompts. Algorithm 1 summarizes how we parse from a prompt to
a belief graph and allow user interaction¹. All agents in §4.1.2 use the same kind of belief graphs.

- Entities. In addition to (a) entities mentioned in
 the user prompt, a belief graph also includes (b)
 implicit entities not mentioned in the prompt but
 likely to appear, e.g., *pet owner* in the context of
 a pet-related scene; and (c) background entities,
 such as *image style, time of day, location*, which
 play important roles in constructing the image.
- Attributes and relations. While the prompt
 might mention some attributes of a certain entity, they are not enough to describe the exact
 details of that entity. Hence the agent have to
 imagine the relevant attributes for each entity,
 and construct a list of possible values along with

269

- 1: **Input:** Initial Prompt (IP)
- 2: **Initialization:** Merged Prompt (MP) \leftarrow IP
- 3: for $turn \leftarrow 1$ to max_turn do
- 4: Parse entities from MP (D.8)
- 5: Parse entity attributes and relations from entities and MP (D.9, D.10)
- 6: Display belief graph, and collect interaction feedback (F)
- 7: Update MP: MP \leftarrow MP + F (D.12)
- 8: **end for**

their associated probabilities (e.g., the *color* attribute for the *cat* entity might have values like *black*, white, gray with corresponding probabilities). Similarly the agent may have to imagine the possible relations between entities, e.g., *spatial relation* between *rabbit* and *cat* might include values like *close*, *far*, *touching*.

¹The clarification question part of the interaction is omitted for simplicity

Importance scores. While the agent can be uncertain about many aspects of the user's intended
image, some are more important than others. E.g., for prompt "a rabbit and a cat", the agent might be
very uncertain about the exact color of a carpet that might appear in the image, but *rabbit* and *cat* are
more important than the carpet. We enable agents to estimate an importance score for each entity,
attribute and relation.

Extracting beliefs and enabling interactions. A simple idea is to use a large language model (LLM) via in-context learning. §D.1 details how an LLM may analyze the user prompt to identify entities, their attributes, and the relations between them, effectively translating the natural language input into a structured representation within the belief. Once the belief is extracted, a user can edit the belief to adjust uncertainty levels, confirm existence of entities etc, as shown in Figure 1.

280 281 4.3 TRANSITION

282 The agent belief undergoes a transition whenever the agent receives new information through user 283 feedback, either user answers from the agent question or user interactions with the graph-based belief interface (Figure 1). This transition process integrates information from the initial user prompt, 284 the conversation history, interaction and the previous belief to generate an updated belief of the 285 user's desired image. Two possible approaches include: (i) Generate a comprehensive prompt that 286 summarizes all interactions and information gathered thus far. This merged prompt is then used to 287 re-generate the belief, effectively incorporating the new information into a refreshed representation. 288 (ii) Leverage natural language to describe the accumulated information, including the initial prompt, 289 conversation history, and user interactions. This descriptive summary is then provided as input to an 290 LLM, instructing it to generate an updated belief based on the provided context. We use (i) for all 291 agents in §4.1.2 and the implementation details can be found in §D.3.

292 293 294

295

296

297

298

5 EXPERIMENTS

We conduct 2 types of experiments to study the effectiveness of the proposed agent design: **automatic evaluation** which uses a simulated user to converse with a T2I agent and **human study** which studies the efficacy of our framework with human subjects.

299 5.1 AUTOMATIC EVALUATION

300 We simulate the user-agent conversation using self-play (Shah et al., 2018) between two LLMs. The 301 conversation starts with an arbitrarily chosen image to represent the goal image from a T2I model 302 that the user has in mind². Along with this ground truth image, a user has a *detailed* prompt in mind 303 that describes the image in high-detail. We use the algorithm similar to Ag2 (detailed in §D.4) to 304 simulate the user, where the questions are answered based on the ground truth prompt and the belief 305 graph generated from the ground truth prompt. We run the agent-user conversation for a total of 15 turns³ and compute different metrics at the end of each turn. More details of the simulated user can 306 be found in the appendix, including the prompts provided to the LLM when simulating the user are 307 provided. Figure 2 part b shows the multi-turn set up that we use in our results. 308

309

310 5.1.1 SETUPS FOR AGENTS AND BASELINE

Baselines. We use a standard T2I model as a baseline, which directly generates an image based on a prompt without asking any questions. We refer to this baseline as 'T2I'.

Agents. We use Ag1, Ag2 and Ag3 with question-asking strategies introduced in §4.1.2. The creation and updates to the belief graph (§4.2), as well as transitions to prompt (§4.3) are consistent among all multi-turn agents. Further implementation details of each agent can be found in §D.

Model Selection. In this work we use an off-the shelve Text-to-Image (T2I) model and a MultiModal Large Language (MLLM) model and build the different components of our agent on top of
these models. We keep these models consistent across all agents for fair comparison. We implement
the agent on top of the Gemini 1.5 (Gemini Team Google, 2024) using the default temperature and a
32K context length. The in-context examples and the exact prompt used at each step of the agent

321

 ²This assumption only applies to the experiments. In practice, users don't necessarily have an image in mind,
 but they can get inspirations from the belief graphs and questions.

³While 15 turns is a suggested approximation of interaction time, accounting for varying difficulty between images, any number of turns can be used with this evaluation approach.



a) generated outputs and target image

b) multi-turn Ag3 example - real generated outputs

Figure 2: a) Each column displays the output of an agent after 15 turns - the right most column shows target image. Target images are part of DesignBench. b) A visualization of the multi-turn set up in the experiments. These are real generated outputs and simulated user outputs at turns 3, 10 and 15.

pipeline is detailed in §D.8 - §D.15. More agent implementation details are provided in §D. For T2I generation, we use Imagen 3 (Baldridge et al., 2024) across all baselines given it's recency and prompt-following capabilities. We used both the models served publically using the Vertex API⁴.

5.1.2 DATASETS.

354 Our multi-turn agents aim to facilitate the generation of complex images, a process that often requires 355 users to iteratively refine text-to-image (T2I) prompts until the generated image aligns with their 356 mental picture. To evaluate these agents, we curate datasets comprising complex scenes involving 357 multiple subjects, interactions, backgrounds, and styles. Each dataset consists of tuples: $(\mathbf{I}, p_0, c, b_{at})$, 358 where I represents the target image, p_0 is an initial (basic) prompt describing only the primary 359 elements of the scene, c is a ground truth caption providing a detailed description of I, including 360 spatial layout, background elements, and style, and b_{at} is the ground truth belief graph constructed via parsing c. The initial prompt p_0 is intentionally less detailed than c to necessitate multi-turn 361 refinement. This framework allows us to assess the agent's ability to guide the user towards the target 362 image I starting from a simplified prompt. 363

364 Existing image-caption datasets primarily focus on simple scenes (Deng et al., 2009; Krizhevsky 365 et al., 2009; Deng, 2012) or focus on very specific categories (Liu et al., 2016; Liao et al., 2022). With 366 the aim for complex realistic images for testing the robustness of the Agents, we evaluate over the validation split of the Coco-Captions dataset (Chen et al., 2015). Five independent human generated 367 captions are provided for each image in the dataset. These captions are often short and describe the 368 basic elements contained in the image and the interactions between objects or persons in the image. 369 We therefore select the shortest of the five human-generated captions and use this as a *starting prompt* 370 p_0 . We then use Gemini 1.5 Pro to expand the starting prompt by adding more details of the attributes 371 of the entities in the image as well as the style and image composition which results in the ground 372 truth caption. We also use the ImageInWords (Garg et al., 2024) dataset which takes a diverse set of 373 realstic and cartoon images and has human annotators create dense detailed captions that describe 374 attribute and relationships between objects in the image. In ImageInWords evaluations we use the 375 long human annotation as the ground truth caption. 376

377

343 344

345

346

347 348

349

350

351 352

⁴https://cloud.google.com/vertex-ai

378 While COCO-Captions and ImageInWords provide complex, real-world images across diverse 379 backgrounds, it lacks the artistic and non-photorealistic imagery often desired by designers and artists 380 seeking to generate content outside the distribution of typical training data. To better evaluate our 381 target for flexible use cases such as by artists, we introduce **DesignBench**, a novel dataset comprising 382 30 scenes specifically designed for this purpose. Each scene follows the (\mathbf{I}, p_0, c) format described earlier. DesignBench includes a mix of cartoon graphics, photorealistic yet improbable scenes, and 383 artistic photographic images. Examples from DesignBench and a comparison with COCO-Captions 384 are provided in the Appendix. 385

386 387 5.1.3 METRICS

The outputs produced by the agent include a final generated image, a final caption and a final belief graph. We evaluate the agents across these modalities and evaluate their alignment to the ground truth image \mathbf{I} , c and b_{qt} , using the following metrics.

Text-Text Similarity: We use 2 metrics for comparing the ground truth caption and the generated caption: 1) T2T – embedding-similarity computed using Gemini 1.5 Pro⁵ and 2) DSG (Cho et al., 2024) adapted to parse text prompts into Davidsonian scene graph using the released code.

Image-Image Similarity (I2I): We compute cosine similarity between the groundtruth image and the generated image from the agent prompt. We use image features from DINOv2 (Oquab et al., 2024) model following prior works.

Text-Image Similarity: We compare the ground truth prompt with the generated image (T2I) using
 the VQAScore (Lin et al., 2024) metric. We use the author released implementation of the metric and
 use Gemini 1.5 Pro as the underlying MLLM. More details about the T2I metrics can be found in §E.

Negative log likelihood (NLL): We construct the ground truth state of the image in the form of a
 belief graph but with no uncertainty. We then approximately compute the NLL of the ground truth
 state given the belief of the agent at each turn, by assuming the independence of all entities, attributes
 and relations, and summing their log probabilities⁶.

405

428

406 5.2 RESULTS FROM AUTOMATED EVALUATION 407

The results from the automatic evaluations in Table 1 show the **I**, c and b_{gt} against each agents final generated image, text and state. All show the mean and standard deviation of the similarity metric at the final agent state. The blue row shows the baseline method which performs no updates to the prompt and instead applies the T2I model to the first prompt. Therefore this baseline represents the lower bound performance.

To add quantitative validity to the ground truth caption generation we perform Text to Image (VQA) Similarity between the ground truth caption and the ground truth over all images in the DesignBench dataset. The mean T2I VQA similarity between the ground truth caption and ground truth image is 0.99999985 with a median 1.0, and standard deviation of 4.5e-07. The mean is extremely close to 1 as expected of an accurate and well formed caption. These numbers can be compared to the T2I column of Table 1 to observe the delta between the ground truth caption and generated captions.

The results in Table 1 show that significant gains in performance come from using proactive multi-turn agents. The blue row shows the simplest baseline which directly uses a T2I model and performs no updates to the initial prompt p_0 . We see that all of the multi-turn agents far exceed the baseline T2I model on both datasets and all metrics. Ag3 (the LLM agent that does not explicitly utilize the belief graph) show superior performance across all metrics.

The plots in Figure 3 show the T2T, I2I, T2I and NLL metrics, averaged across all images in the ImageInWords dataset, per turn for 15 turns. We see that the multi-turn agents all improve in every metric as they increase the number of interactions. Interestingly we see the T2T and the T2I VQA similarity metric seems to plateau or decrease after 10 interactions, while the I2I scores continue to

⁵Text embeddings are obtained from Embeddings API: https://ai.google.dev/gemini-api/docs/embeddings.

 ⁶This approximation does not account for potential similarities in the names of entities or attributes. This could lead to approximation errors if, for example, the model confuses "Persian cat" with "Siamese cat" due to their similar names. Addressing this limitation would require incorporating semantic similarity measures into the NLL computation.

Dataset	Model	T2T \uparrow	I2I (DINO) \uparrow	T2I (VQAScore)↑	NLL↓	DSG (T2T)↑
	T2I	$0.8757{\pm}.03$	$0.5170 {\pm}.16$	$0.2976 {\pm}.45$	520.0645 ± 161.3	$0.5904 {\pm} .05$
Coco Contions	Ag1	$0.9440{\pm}.02$	$0.6269 {\pm}.12$	$0.5831 {\pm}.49$	$508.4014{\pm}158.5$	$0.7555{\pm}.08$
Coco-Captions	Ag2	$0.9461{\pm}.02$	$0.6141{\pm}.13$	$0.6632 {\pm}.46$	$481.7224{\pm}154.5$	$0.8344{\pm}.08$
	Ag3	$\textbf{0.9501} {\pm}.02$	$0.6575 \pm .10$	0.7751 ±.39	$\textbf{446.5679}{\pm}151.8$	$\textbf{0.9001} {\pm}.05$
	T2I	$0.8807 {\pm} .02$	$0.5154 {\pm}.15$	$0.3711 {\pm}.47$	$459.9053{\pm}200.2$	$0.6815 {\pm}.70$
ImageInWords	Ag1	$0.9429 {\pm}.02$	$0.5548 {\pm} .15$	$0.5058{\pm}.48$	$449.8927{\pm}196.1$	$0.8162 {\pm} .08$
	Ag2	$0.9382{\pm}.02$	$0.5645 {\pm} .15$	$0.5701 {\pm}.48$	$444.5227 {\pm} 193.7$	$0.8791{\pm}.07$
	Ag3	$\textbf{0.9418} {\pm}.02$	$0.5875 \pm .14$	$0.6624 \pm .45$	$\textbf{429.4636}{\pm}194.5$	$\textbf{0.9124} {\pm}.06$
DesignBench	T2I	$0.8740 {\pm}.02$	$0.5439 {\pm}.12$	$0.3528 {\pm}.48$	$320.8898 {\pm} 93.7$	$0.6074 {\pm} .08$
	Ag1	$0.9365{\pm}.02$	$0.5943{\pm}.12$	$0.6848 {\pm}.46$	$295.1974{\pm}69.2$	$0.8285{\pm}.08$
	Ag2	$0.9384{\pm}.02$	$0.6417 {\pm}.11$	$0.8553 {\pm} .34$	271.2604 ± 81.9	$0.9181{\pm}.06$
	Ag3	$\textbf{0.9429} {\pm}.02$	0.6924 ±.12	0.9545 ±.21	257.4352 ± 67.5	0.9485 ±.04

increase. The NLL metric shows large performance gains of the Ag3 agent in comparison to all other methods. The plots in Figure 10 shows the T2T DSG metrics.

Table 1: Automatic evaluation results on Coco-Captions, ImageInWords, and DesignBench. Agents show large performance gains in all metrics over a standard T2I model alone.

5.3 ANALYSIS OF QUANTITATIVE RESULTS

454 The evaluations on the 455 COCO-captions, ImageIn-456 Words, DesignBench 457 datasets similar show 458 results and highlight the 459 same patterns across the 460 different agents.

461 **Multi-Turn Agents show** 462 clear advantage: The im-463 mediate take away is the 464 baseline which does not use 465 multi-turn interaction and 466 instead passes in the original prompt into the T2I 467 model performs worse than 468 the multi-turn agents on all 469 metrics on both datasets. 470 This confirms our hypoth-471 esis that the current T2I 472 agents often produce less 473 desirable images given am-474





biguity in prompts. In Figure 2 we see real outputs of the multi-turn set up with the Ag3 agent.

LLMs being a part of agents play a significant role: The best performers (Ag2 and Ag3) both query 476 and LLM to provide a question to ask the user based on contextual information such as the belief 477 graph and conversation history. They query the LLM to construct a concise and clear question but 478 don't impose further constraints on the question construction. Ag1 provides a programatic template 479 for how the LLM should construct the question based on its belief graph and does not provide any 480 conversation history information. Examples of dialogs and the generated questions produced by the 481 three agents can be found in the Appendix in Figure 4. This figure demonstrates that the templated 482 question creation leads to extremely specific questions that often gather minimal information in return. 483 This is an intrinsic limitation of hard coded question selection strategy but also can be an issue of the heuristic scores we defined for question selection in Ag1. In contrast, Ag2 and Ag3 generate 484 questions that are more open-ended thus allowing the user to provide more nuanced details which in 485 consequence enhance the Agent's image knowledge.

434

449

450

451 452

453

Feature	V. Unlikely (%)	Unlikely (%)	Could Help (%)	Likely (%)	V. Likely (%)
Clarifications	3.5	5.6	31.5	37.8	21.7
Entity Graph	4.2	7.7	35	32.9	20.3
Relation Graph	7	7	37.1	28.7	20.3

Table 2: Perceived helpfulness of proposed features (% of users) rated by 143 raters.

Question prompts with question-asking principles show advantage over those with beliefs: The 494 Ag3 agent (which uses an LLM with question generation instructions about entity, attributes etc related to the belief) dominates across both datasets on every metric. Ag2 uses the belief explicitly to 495 construct questions by passing the belief into the LLM as information from which to generate the next question. When inspecting the reasoning steps of Ag2, we found that Ag2 excessively relies on importance scores in beliefs to ask questions, and if the importance scores are not estimated properly, 498 the quality of the questions decreases. 499

5.4 HUMAN STUDY

In order to get real user feedback, we performed a human survey with the objective of understanding user frustrations to validate whether our potential solutions could help with their use of T2I models. 504 We gathered data from 143 participants who all identified to be regular T2I users (at least once a 505 month). Participants were presented with four hypothesized frustrations (prompt misinterpretation, 506 many iterations, inconsistent generations, incorrect assumptions) and three potential mitigating 507 features (clarifications, entity graph, relationship graph; more details in Appendix §H). 508

As reported in Table 4 (in Appendix), the results confirmed the prevalence of hypothesized frustrations 509 amongst users, with 83% experiencing occasional, frequent, or very frequent frustration due to 510 prompt iterations, followed by 70% for misinterpretations, 71% for inconsistent generations, and 60% 511 experiencing frustration due to incorrect assumptions. Most acutely 55% of participants reported 512 frequent or very frequent frustration due to the prompt iteration frequency necessary. In Table 2, we 513 report the mitigation features that are likely to help. Clarifications reported the highest likelihood 514 to help current workflows (91% could / likely / very likely to be helpful), followed by entity graphs 515 (88% could / likely / very likely to be helpful) and relationship graphs (86% could / likely / very 516 likely to be helpful). Clarifications were expected to deliver value immediately / very soon by 58%.

517 Overall these suggest strong user desire for & likelihood for success of features that reduce iterations 518 and mitigate misinterpretations in T2I generation. Full explanations of the hypothesized frustrations, 519 mitigation and responses splits are in §H. All respondents were compensated for their time as per 520 market rates, and were recruited by our vendor to ensure diversity across age, gender, and T2I usage 521 in terms of models, frequency and purpose (work and non work).

522 523 524

525

526

527

528

492 493

496

497

500

501 502

6 DISCUSSION AND CONCLUSION

This work introduces a design for agents that assist users in generating images through an interactive process of question-asking and belief graph refinement. By dynamically updating its understanding of the user's intent, the agent facilitates a more collaborative and precise approach to image generation.

529 Modular design. Our agent prototypes are highly modular: the agents use frozen T2I models to 530 generate images based on the prompts that the agent updated. Therefore when a better off-the-shelf T2I model becomes available, it can be directly plugged into the agents and the system will achieve 531 better performance without any additional adaptation⁷. 532

533 Future work. Alternative to the modular design, one can explore generating images directly from 534 belief graphs and fine-tuning LLM/VLMs on text/image trajectories that include asking questions. These may require a) collecting data such as gold-standard trajectories or annotations on the quality of 536 trajectories of human-agent conversations and b) new approaches to fine-tune the model on multi-turn 537 trajectories of images and text, which can potentially improve the performance of the agent.

⁷T2T scores in Table 1 ablates the T2I model and only performs similarity on the captions. Our agents have achieved a 92%+ T2T score, showing that their performance can be boosted by adopting better T2I models.

540 ETHICS STATEMENT 541

542 Our proposed T2I agents are equipped with better tools (belief graphs) for interpretability and 543 controllability. Presenting the agent's belief graph can be a generalizable method for AI transparancy, 544 which is an important factor given the increasing complexity of modern AI models.

⁵⁴⁵ By asking clarification questions, our proposed agents may enable a more customizable and person⁵⁴⁶ alized content creation experience. Because different groups of people may perceive harmfulness
⁵⁴⁷ of contents differently, learning more about the user through clarification questions can potentially
⁵⁴⁸ mitigate risks of generating contents that can be offensive to each specific user.

550 REPRODUCIBILITY

549

561

562

567

568

569

570

574

575

576

577

We plan to release all code and DesignBench upon publication. All implementation details and
prompts used in this work can be found in the appendix. All models we used in this work are publicly
accessible with APIs linked in the experiments.

555 ACKNOWLEDGEMENTS

We would like to thank Jason Baldridge and Zoubin Ghahramani for insightful discussions on multi turn T2I and belief states, Mahima Pushkarna for the help and consultation on user study. We would also like to thank Richard Song and Noah Fiedel for feedback on the paper.

References

- Yusra Alkhazraji, Matthias Frorath, Markus Grützner, Malte Helmert, Thomas Liebetraut, Robert Mattmüller, Manuela Ortlieb, Jendrik Seipp, Tobias Springenberg, Philip Stahl, and Jan Wülfing.
 Pyperplan. https://doi.org/10.5281/zenodo.3700819, 2020. URL https:// doi.org/10.5281/zenodo.3700819.
 - Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N Bennett, Kori Inkpen, et al. Guidelines for human-ai interaction. In ACM Conference on Human Factors in Computing Systems (CHI), 2019.
- Jason Baldridge, Jakob Bauer, Mukul Bhutani, Nicole Brichtova, Andrew Bunner, Kelvin Chan,
 Yichang Chen, Sander Dieleman, Yuqing Du, Zach Eaton-Rosen, et al. Imagen 3. *arXiv preprint arXiv:2408.07009*, 2024.
 - James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3):8, 2023.
- Carrie J Cai, Emily Reif, Narayan Hegde, Jason Hipp, Been Kim, Daniel Smilkov, Martin Wattenberg,
 Fernanda Viegas, Greg S Corrado, Martin C Stumpe, et al. Human-centered tools for coping with
 imperfect algorithms during medical decision-making. In ACM Conference on Human Factors in
 Computing Systems (CHI), 2019.
- 582
 583
 584
 Kathryn Chaloner and Isabella Verdinelli. Bayesian experimental design: A review. *Statistical science*, pp. 273–304, 1995.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- Yanda Chen, Ruiqi Zhong, Narutatsu Ri, Chen Zhao, He He, Jacob Steinhardt, Zhou Yu, and Kathleen McKeown. Do models explain themselves? counterfactual simulatability of natural language explanations, 2023. URL https://arxiv.org/abs/2307.08678.
- Yida Chen, Aoyu Wu, Trevor DePodesta, Catherine Yeh, Kenneth Li, Nicholas Castillo Marin, Oam
 Patel, Jan Riecke, Shivam Raval, Olivia Seow, et al. Designing a dashboard for transparency and control of conversational ai. *arXiv preprint arXiv:2406.07882*, 2024.

594 595 596	Junhao Cheng, Xi Lu, Hanhui Li, Khun Loun Zai, Baiqiao Yin, Yuhao Cheng, Yiqiang Yan, and Xiaodan Liang. Autostudio: Crafting consistent subjects in multi-turn interactive image generation. <i>arXiv preprint arXiv:2406.01388</i> , 2024a.
597 598 599	Junhao Cheng, Baiqiao Yin, Kaixin Cai, Minbin Huang, Hanhui Li, Yuxin He, Xi Lu, Yue Li, Yifei Li, Yuhao Cheng, et al. Theatergen: Character management with llm for consistent multi-turn image generation. arXiv preprint arXiv:2404.18919, 2024b.
600 601 602	Jaemin Cho, Yushi Hu, Roopal Garg, Peter Anderson, Ranjay Krishna, Jason Baldridge, Mohit Bansal,
603 604	evaluation for text-image generation. <i>arXiv preprint arXiv:2310.18235</i> , 2023.
605 606 607 608	Jaemin Cho, Yushi Hu, Roopal Garg, Peter Anderson, Ranjay Krishna, Jason Baldridge, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. Davidsonian scene graph: Improving reliability in fine- grained evaluation for text-to-image generation, 2024. URL https://arxiv.org/abs/ 2310.18235.
609 610 611 612 613 614	Siddhartha Datta, Alexander Ku, Deepak Ramachandran, and Peter Anderson. Prompt expansion for adaptive text-to-image generation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 3449–3476, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.acl-long.189.
615 616 617	Donald Davidson. Theories of meaning and learnable languages. In Yehoshua Bar-Hillel (ed.), <i>Proceedings of the 1964 International Congress for Logic, Methodology, and Philosophy of Science</i> , pp. 383–394. North-Holland Publishing, 1965.
618 619 620	Donald Davidson. Truth and meaning. In <i>Philosophy, Language, and Artificial Intelligence: Resources for Processing Natural Language</i> , pp. 93–111. Springer, 1967a.
621 622	Donald Davidson. The logical form of action sentences. The Logic of Decision and Action, 1967b.
623 624 625	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In <i>IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 248–255. IEEE, 2009.
626 627	Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. <i>IEEE signal processing magazine</i> , 29(6):141–142, 2012.
628 629 630	Richard E Fikes and Nils J Nilsson. Strips: A new approach to the application of theorem proving to problem solving. <i>Artificial intelligence</i> , 2(3-4):189–208, 1971.
631 632 633 634	Roopal Garg, Andrea Burns, Burcu Karagol Ayan, Yonatan Bitton, Ceslee Montgomery, Yasumasa Onoe, Andrew Bunner, Ranjay Krishna, Jason Baldridge, and Radu Soricut. Imageinwords: Unlocking hyper-detailed image descriptions, 2024. URL https://arxiv.org/abs/2405. 02793.
635 636	Roman Garnett. Bayesian optimization. Cambridge University Press, 2023.
637 638 639 640	Caelan Reed Garrett, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. PDDLstream: Integrating symbolic planners and blackbox samplers via optimistic adaptive planning. In <i>International Conference on Automated Planning and Scheduling</i> , volume 30, pp. 440–448, 2020a.
641 642 643	Caelan Reed Garrett, Chris Paxton, Tomás Lozano-Pérez, Leslie Pack Kaelbling, and Dieter Fox. Online replanning in belief space for partially observable task and motion problems. In <i>IEEE International Conference on Robotics and Automation (ICRA)</i> , 2020b.
644 645	Gemini Team Google. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. URL https://arxiv.org/abs/2403.05530.
647	Matthew L Ginsberg and David E Smith. Reasoning about action i: A possible worlds approach. <i>Artificial Intelligence</i> , 35(2):165–195, 1988.

648 649 650	Naresh Gupta and Dana S Nau. On the complexity of blocks-world planning. <i>Artificial Intelligence</i> , 56(2-3):223–254, 1992.
651 652 653	Melissa Hall, Candace Ross, Adina Williams, Nicolas Carion, Michal Drozdzal, and Adriana Romero Soriano. Dig in: Evaluating disparities in image generations with indicators for geographic diversity, 2024. URL https://arxiv.org/abs/2308.06198.
654 655	Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference- free evaluation metric for image captioning, 2022.
655 657 658	Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018.
659 660	Kristina Höök. Steps to take before intelligent user interfaces become real. <i>Interacting with computers</i> , 12(4):409–426, 2000.
661 662 663 664	Hexiang Hu, Kelvin CK Chan, Yu-Chuan Su, Wenhu Chen, Yandong Li, Kihyuk Sohn, Yang Zhao, Xue Ben, Boqing Gong, William Cohen, et al. Instruct-imagen: Image generation with multi-modal instruction. In <i>International Conference on Computer Vision (ICCV)</i> , 2024.
665 666 667	Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. TIFA: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In <i>International Conference on Computer Vision (ICCV)</i> , 2023.
668 669 670	Minbin Huang, Yanxin Long, Xinchi Deng, Ruihang Chu, Jiangfeng Xiong, Xiaodan Liang, Hong Cheng, Qinglin Lu, and Wei Liu. Dialoggen: Multi-modal interactive dialogue system for multi-turn text-to-image generation. <i>arXiv preprint arXiv:2403.08857</i> , 2024.
671 672 673 674 675	Ben Hutchinson, Jason Baldridge, and Vinodkumar Prabhakaran. Underspecification in scene description-to-depiction tasks. In <i>Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , 2022.
676 677	Leslie Pack Kaelbling and Tomás Lozano-Pérez. Hierarchical task and motion planning in the now. In <i>IEEE International Conference on Robotics and Automation (ICRA)</i> , 2011.
678 679	Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. Journal or Artificial Intelligence Research (JAIR), 4:237–285, 1996.
681 682	Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in partially observable stochastic domains. <i>Artificial Intelligence</i> , 101(1-2):99–134, 1998.
683 684 685	Nithish Kannen, Arif Ahmad, Marco Andreetto, Vinodkumar Prabhakaran, Utsav Prabhu, Adji Bousso Dieng, Pushpak Bhattacharyya, and Shachi Dave. Beyond aesthetics: Cultural competence in text-to-image models, 2024. URL https://arxiv.org/abs/2407.06863.
687 688 689	Sunnie SY Kim, Elizabeth Anne Watkins, Olga Russakovsky, Ruth Fong, and Andrés Monroy- Hernández. "help me help the ai": Understanding how explainability can support human-ai interaction. In <i>ACM Conference on Human Factors in Computing Systems (CHI)</i> , 2023.
690 691	Roger E Kirk. Experimental design. <i>The SAGE Handbook of Quantitative Methods in Psychology</i> , pp. 23–45, 2009.
692 693	Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a- pic: An open dataset of user preferences for text-to-image generation, 2023.
695	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
696 697 698 699	Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models, 2019. URL https://arxiv.org/abs/1904.06991.
700 701	Zeqiang Lai, Xizhou Zhu, Jifeng Dai, Yu Qiao, and Wenhai Wang. Mini-dalle3: Interactive text to image by prompting large language models, 2023. URL https://arxiv.org/abs/2310.07653.

702 703 704 705 706 707 708	Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Her- nandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamilė Lukošiūtė, Karina Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson, Sam McCandlish, Sandipan Kundu, Saurav Kadavath, Shannon Yang, Thomas Henighan, Timothy Maxwell, Timothy Telleen-Lawton, Tristan Hume, Zac Hatfield-Dodds, Jared Kaplan, Jan Brauner, Samuel R. Bowman, and Ethan Perez. Measuring faithfulness in chain-of-thought reasoning, 2023. URL https://arxiv.org/abs/2307.13702.
709 710 711	Peiyuan Liao, Xiuyu Li, Xihui Liu, and Kurt Keutzer. The artbench dataset: Benchmarking generative models with artworks. <i>arXiv preprint arXiv:2206.11404</i> , 2022.
712 713 714	Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In <i>Europ. Conference on Computer Vision (ECCV)</i> , 2014.
715 716 717 718	Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation, 2024. URL https://arxiv.org/abs/2404.01291.
719 720 721	Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 1096–1104, 2016.
722 723 724	Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and William Yang Wang. Llmscore: Unveiling the power of large language models in text-to-image synthesis evaluation, 2023.
725 726 727	John McCarthy and Patrick J. Hayes. Some philosophical problems from the standpoint of artificial intelligence. In B. Meltzer and D. Michie (eds.), <i>Machine Intelligence 4</i> , pp. 463–502. Edinburgh University Press, 1969. reprinted in McC90.
728 729 730	Marvin Minsky. A framework for representing knowledge. Technical report, A.I. Laboratory, Mas- sachusetts Institute of Technology, 1974. URL https://dspace.mit.edu/bitstream/ handle/1721.1/6089/AIM-306.pdf.
731 732	Marvin Minsky. The Society of Mind. Simon and Schuster, 1988.
733 734 735	Robert C Moore. A formal theory of knowledge and action. <i>Formal theories of the commonsense world</i> , pp. 319–358, 1985.
736 737 738	Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models, 2020. URL https://arxiv.org/abs/2002.09797.
739 740	Nils J Nilsson. Probabilistic logic. Artificial Intelligence, 28(1):71-87, 1986.
741	Nils J Nilsson. The Quest for Artificial Intelligence. Cambridge University Press, 2009.
742 743 744	Donald A Norman. How might people interact with agents. <i>Communications of the ACM</i> , 37(7): 68–71, 1994.
745 746 747 748 749 750	Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. URL https://arxiv.org/abs/2304.07193.
751 752 753	Hanna M Pasula, Luke S Zettlemoyer, and Leslie Pack Kaelbling. Learning symbolic models of stochastic domains. <i>Journal or Artificial Intelligence Research (JAIR)</i> , pp. 309–352, 2007.
754 755	Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. <i>arXiv preprint arXiv:2307.01952</i> , 2023.

756	Stuart J Russell and Peter Norvig. Artificial intelligence: A modern approach. Pearson, 2016.
758 759	Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans, 2016.
760 761	Burr Settles. Active learning literature survey. Computer Sciences Technical Report 1648, University of Wisconsin–Madison, 2009.
762 763 764 765	Pararth Shah, Dilek Hakkani-Tür, Gokhan Tür, Abhinav Rastogi, Ankur Bapna, Neha Nayak, and Larry Heck. Building a conversational agent overnight with dialogue self-play, 2018. URL https://arxiv.org/abs/1801.04871.
766 767	Heyu Sun and Qiang Guo. Dsg-gan: Multi-turn text-to-image synthesis via dual semantic-stream guidance with global and local linguistics. <i>Intelligent Systems with Applications</i> , 20:200271, 2023.
768 769	Richard S Sutton. Reinforcement learning: An introduction. A Bradford Book, 2018.
770 771	Fernanda Viégas and Martin Wattenberg. The system model and the user model: Exploring ai dashboard design. <i>arXiv preprint arXiv:2305.02469</i> , 2023.
772 773 774	Kailas Vodrahalli and James Zou. Artwhisperer: A dataset for characterizing human-ai interactions in artistic creations. In <i>International Conference on Machine Learning (ICML)</i> , 2024.
775 776 777	Qianli Wang, Tatiana Anikina, Nils Feldhus, Josef van Genabith, Leonhard Hennig, and Sebastian Möller. Llmcheckup: Conversational examination of large language models via interpretability tools and self-explanations, 2024a. URL https://arxiv.org/abs/2401.12576.
778 779 780	Zekang Wang, Li Liu, Huaxiang Zhang, Dongmei Liu, and Yu Song. Generative adversarial text-to- image generation with style image constraint. <i>Multimedia Systems</i> , 29(6):3291–3303, 2023.
781 782 783	Zi Wang, Alexander Ku, Jason Baldridge, Tom Griffiths, and Been Kim. Gaussian process probes (GPP) for uncertainty-aware probing. In <i>Advances in Neural Information Processing Systems</i> (<i>NeurIPS</i>), 2024b.
784 785 786	Fanyue Wei, Wei Zeng, Zhenyang Li, Dawei Yin, Lixin Duan, and Wen Li. Powerful and flexible: Per- sonalized text-to-image generation via reinforcement learning. <i>arXiv preprint arXiv:2407.06642</i> , 2024.
787 788 789 790	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023a. URL https://arxiv.org/abs/2201.11903.
791 792 793	Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning differently, 2023b. URL https://arxiv.org/abs/2303.03846.
794 795 796 797	Olivia Wiles, Chuhan Zhang, Isabela Albuquerque, Ivana Kajić, Su Wang, Emanuele Bugliarello, Yasumasa Onoe, Chris Knutsen, Cyrus Rashtchian, Jordi Pont-Tuset, et al. Revisiting text-to-image evaluation with gecko: On metrics, prompts, and human ratings. <i>arXiv preprint arXiv:2404.16820</i> , 2024.
798 799 800 801	Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis, 2023.
802 803 804	Victoria Xia, Zi Wang, Kelsey Allen, Tom Silver, and Leslie Pack Kaelbling. Learning sparse relational transition models. In <i>International Conference on Learning Representations (ICLR)</i> , 2019.
805 806 807	Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation, 2023.
808 809	Qian Yang, Aaron Steinfeld, Carolyn Rosé, and John Zimmerman. Re-examining whether, why, and how human-ai interaction is uniquely difficult to design. In <i>ACM Conference on Human Factors in Computing Systems (CHI)</i> , 2020.

Xi Ye and Greg Durrett. The unreliability of explanations in few-shot prompting for textual reasoning, 2022. URL https://arxiv.org/abs/2205.03401.

Guojun Yin, Bin Liu, Lu Sheng, Nenghai Yu, Xiaogang Wang, and Jing Shao. Semantics disentangling
 for text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2327–2336, 2019.

Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian Karrer, Shelly Sheynin, et al. Scaling autoregressive multi-modal models: Pretraining and instruction tuning. *arXiv preprint arXiv:2309.02591*, 2(3), 2023.

Lai Zeqiang, Zhu Xizhou, Dai Jifeng, Qiao Yu, and Wang Wenhai. Mini-dalle3: Interactive text to image by prompting large language models. *arXiv preprint arXiv:2310.07653*, 2023.

Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. Alignscore: Evaluating factual consistency with a unified alignment function, 2023.

864 A NOVELTY AND CONTRIBUTIONS

866 867	In this section, we emphasize the novelty and contributions of this work.
868	1 System design of presentive T2I agents
869	1. System design of proactive 121 agents.
870	• Novel numan-agent interaction modalities: Prior to our work, numan users typically interact with current T2L systems by giving additional instructions or refining the
871	prompt. To the best of our knowledge, our work is the first to propose a proactive T2I
872	agent system that is able to ask clarification questions and present its belief graph for
873	the user to edit.
874	• Novel human-agent interaction interface: We designed a new interface to best enable
875	the clarification and belief graph interaction modalities. We have not seen these features
876	in any T2I, or other generative media apps that are publicly live to date, signifying to
877	us total uniqueness. Our human studies showed that at least 85
878	• Novel design of different T2I agents that enable the proposed interaction modalities.
879	Please see Section 4 of the paper for the full details of the design principles and
880	construction of those T2I agent prototypes (Ag1, Ag2, Ag3).
881	2. Our belief graph significantly differs from the classic belief state in the following ways:
882	• Hardcoded predicates v.s. Automatically-generated predicates: Traditionally, con-
883	structing classic symbolic belief states requires a pre-defined set of predicates such
884	as "on(a, b)", "is_red(a)", "at_position(robot, x, y, z)" and it is non-trivial to learn new
885	predicates that can be used and generalized to new tasks (Pasula et al., 2007; Xia et al.,
886	2019). Typically the pre-defined set of predicates are written by system developers and
887	hardcoded into classic AI systems (Fikes & Nilsson, 1971).
888	• Our benefigraph does not require any pre-defined set of predicates. Instead, we propose to construct symbolic beliefs using a sequential in context learning (ICL) method with
889	I I Ms. This method first generates a list of entities together with their descriptions
890	conditioned on the user description of image: then, we add each entity to the context.
891	and let the LLM generate a list of attributes and values (this step is done in parallel
892	across entities); and finally, we add all entities to the context and let the LLM generate
893	relations and their attributes. Our method can be generalized across a wide range
894	of T2I tasks and achieve high performance (see our comprehensive results on Coco,
895	Imageinwords, DesignBench). We have included a pseudo code for this method in the
896	paper.
897	• Application to planning v.s. 121: To the best of our knowledge, classic symbolic belief
898	states are mostly used for robol planning, and we are the first to use symbolic beliefs to assist T21 tasks. Data structure for symbolic states / balief states - Sat y s. Graph:
899	Because of the application to planning, a symbolic world state is usually implemented
900	and stored as a set or list of literals (atoms or negation of atoms where atoms are
901	instantiated predicates (Alkhazraji et al., 2020; Garrett et al., 2020a;b) so that whenever
902	an action is applied, the agent can apply transition by adding and deleting items in the
903	set according to the precondition and effect of the action.
904	• For T2I tasks, it is more convenient to use a graph to represent the world state associated
905	with an image, since entities and relations naturally form a set of nodes (entities)
906	and edges (relations between entities). Each component of the graph can also have
907	probabilities, making it easy to turn a world state into a belief state using the same data
908	structure. Hence we represent 121 agent beliefs using graphs. The agent can directly
909	• Interpretability and controllability: The graph structure makes our agent belief more
910	interpretability and controllability. The graph structure makes our agent benef more interpretable than traditional belief states, since we can visualize and progressively
911	disclose the graph to the human user. Moreover, each node or relation in the belief graph
912	has associated descriptions, making it easy for the user to understand and potentially
913	edit every component of the belief graph. In our human studies, about 85% of raters
914	found the belief graph useful. To the best of our knowledge, our work is the first to use
915	the graph-based belief state for human-AI interaction.
916	3. Automated evaluation of T2I agents: We propose a novel automated evaluation approach for
917	T2I agents using self-play. The agent interacts with a simulated user that has access to the original image and its long caption. See Section 5.1 (and C.4) for the full details of how the

918	simulated user is constructed. This evaluation pipeline is easy to use and can help the future
919	development of T2I agents.
920	4. DesignBench: We envision that a significant fraction of T2I users are artists and designers,
921	and it is important to ensure that T2I agents are evaluated for these use cases. Hence we
922	create DesignBench, featuring photo-realism, animation and multiple styles with short
923	and long captions. DesignBench can be directly plugged into our automated evaluation to
924	streamline the evaluation process.
925	
920	B FORMALISM OF THE AGENT AND ITS OBJECTIVE
927	
920 929	We define an interactive T2I agent as a $\langle B, A, O, \tau, \pi \rangle$ tuple, where we have
930	• S: a representation space of images,
931	• B: a space of agent beliefs,
932	• A: a space of actions that the agent can take,
933	• O: a space of agent observations of the user,
934	• transition function $\tau: B \times A \times O \mapsto B$ for updating beliefs given new interactions,
935	• action selection strategy $\pi: B \mapsto A$, which specifies which action to take given a belief.
936	For each user-initiated interaction, we assume that there exists a specific intent $s \in S$, where S is
937	the space of all possible user intents. For a T2I task, we assume that the intent is the image the user
938	would like to generate, and the intent stays the same throughout the interaction with an agent. We
939	discuss more about the validity of this assumption in §6.
940	Each type of T2I agents can have a unique user intent representation, belief representation, construc-
941	tion of the action space, and user interface design to obtain observations of users.
942	In §4, we show the examples for these components.
943	
944	We use a score function, $f: B \times S \mapsto \mathbb{R}$, to evaluate the alignment between an agent belief and a user intent at any turn of the interaction. Function f can only be evaluated in hindsight once the user
945	intent is revealed. The agent does not have direct access to function f since the user intent is hidden
0/17	from the agent. However, the agent may construct a probabilistic distribution over function f based
0/18	on its belief about the user intent. The goal of the agent is to maximize function f with as few turns
940	of interaction with the user as possible.
950	
951	C VISUALIZATION OF MULTI-TURN AGENT-USER DIALOGS AND
952	Generated Images
953	OENERALED IMAOES
954	In Figure 4, we show examples of multi-turn dialogs between simulated users and the three agents in
955	Section 5. We also visualize the generated images in Figure 5. Figure 6. Figure 7 and Figure 8
956	section 5. The also visualize the generated images in Figure 5, Figure 6, Figure 7 and Figure 6.
957	

D IMPLEMENTATION DETAILS (FOR ALL AGENTS IN OUR EXPERIMENT)

We propose three distinct T2I agents, each characterized by a unique configuration of $\langle B, A, O, \tau, \pi \rangle$ tuples:

- *Ag1: Heuristic Score Agent:* this agent incorporates a human-defined heuristic score based on the belief to guide question generation. This heuristic score reflects the perceived importance of different aspects of the belief in driving the conversation forward;
 - Ag2: Belief-prompted Agent: This agent leverages an LLM to generate questions by processing both the conversation history and a structured representation of the belief.
- Ag3: Principle-prompted Agent: This agent generates questions directly from the conversation history, relying solely on the implicit knowledge and reasoning capabilities of the underlying Large Language Model (LLM). It does not employ an explicit, structured belief representation;





Figure 5: Agent Generated Image Outputs on DesignBench: a chart of the generated image outputs of the four main Agent types in comparison to the goal image. Each column displays the output of a different agent and the right most column shows the goal image that the agents aimed to recreate. Each agent was provided with the same starting prompt and iterated for 15 turns, with the exception of the "T2I" agent column which produces an image from the starting prompt. Ag1, Ag2 and Ag3 refer to the Agents described in §D. Each agent uses the same T2I model to produce the final image. The goal images displayed here are from our DesignBench dataset described in the experiments section.

- 1078
- 1079



Figure 6: Agent Generated Image Outputs on DesignBench (Continued): a chart of the generated image outputs of the four main Agent types in comparison to the goal image. Each column displays the output of a different agent and the right most column shows the goal image that the agents aimed to recreate. Each agent was provided with the same starting prompt and iterated for 15 turns, with the exception of the "T2I" agent column which produces an image from the starting prompt. Ag1, Ag2 and Ag3 refer to the Agents described in §D. Each agent uses the same T2I model to produce the final image. The goal images displayed here are from the DesignBench dataset described in the experiments section.



Figure 7: Agent Generated Image Outputs (Coco-Captions Validation): a chart of the generated image outputs of the four main Agent types in comparison to the goal image. Each column displays the output of a different agent and the right most column shows the goal image that the agents aimed to recreate. Each agent was provided with the same starting prompt and iterated for 15 turns, with the exception of the "T2I" agent column which produces an image from the starting prompt. Ag1, Ag2 and Ag3 refer to the Agents described in §D. Each agent uses the same T2I model to produce the final image. The goal images displayed here are from the Coco-Captions Chen et al. (2015) dataset described in the experiments section.

- 1185
- 1186
- 1187



Figure 8: Agent Generated Image Outputs (Coco-Captions Validation): a chart of the generated image outputs of the four main Agent types in comparison to the goal image. Each column displays the output of a different agent and the right most column shows the goal image that the agents aimed to recreate. Each agent was provided with the same starting prompt and iterated for 15 turns, with the exception of the "T2I" agent column which produces an image from the starting prompt. Ag1, Ag2 and Ag3 refer to the Agent described in the methods section. Each agent uses the same T2I model to produce the final image. The goal images displayed here are from the Coco-Captions Chen et al. (2015) dataset described in the experiments section.

1242 D.1 IMPLEMENTATION OF STATE

Our agent's state is represented in two complementary forms: (i) Merged prompt: This is a natural language representation that summarizes the entire conversation history up to the current turn. It provides a comprehensive textual overview of the user's requests, feedback, and any clarifications exchanged with the agent. (ii) Belief: This is a symbolic representation derived from the merged prompt. It parses the natural language text into a structured format, capturing key elements like entities, attributes, relationships, and associated probabilities. This structured representation facilitates more precise reasoning and decision-making by the agent.

Prompt Merging. An LLM (§D.11) summarizes the latest interaction, encapsulating the agent's question and the user's response into a concise textual representation. This step distills the essential information exchanged during the interaction. Another LLM (§D.12) merges the summarized interaction with the existing merged prompt, which contains the accumulated information from previous interactions. This creates an updated prompt that reflects the evolving understanding of the user's intent.

1256 Belief Parsing. See an example of the belief state fig. 9. We employ three specialized parsers 1257 trained via in-context learning (ICL): entity parser (§D.8) analyzes the user prompt to identify and 1258 extract a list of relevant entities.; attribute parser (§D.9) takes user prompt and an entity as the input to 1259 extract a list of attributes associated with that entity; relation parser (§D.10) takes the user prompt and 1260 a list of entities as input and identifies relationships between those entities. Each entity is associated 1261 with meta information like name, importance to ask score, description, probability of appearing, a list of attributes like color, position, etc⁸. Each attribute contains meta information like name, importance 1262 to ask score, a list of possible values for the attribute along with their associated probabilities, etc. 1263 Each relation includes meta information such as: name, description, spatial relation, importance to 1264 ask score, entity 1 and entity 2, whether the relation is bidirectional, etc. 1265

- 1266
- 1267
- 1268
- ._03
- 1270 1271
- 1272
- 1273

1274

12/5

1276 1277

1278

- 1280 1281
- 1282
- 1283
- 1284 1285
- 1286
- 1287
- 1288
- 1209
- 1291

 ⁸Name is a unique identifier for the entity; Importance to ask score: A numerical value indicating the entity's perceived importance in satisfying the user's request. Entities with higher scores are prioritized during question generation, as they are likely to reduce uncertainty and contribute significantly to the final image;
 Description provides a textual description of the entity; probability of appearing estimates likelihood of the entity being present in the generated image; Attributes is for understanding the detailed attributes of the entities.

1297	
1298	
1299	
1300	
1301	
1302	
1303	
1304	
1305	
1306	
1307	
1308	Belief state
1309	
1310	Entities
1311	Pabbit
1312	Attribute Name: color, importance, score: 0.9
1313	candidates: [brown: 0.25, white: 0.25, grey: 0.2, black: 0.15,]
1314	Attribute Name: breed, importance_score: 0.3,
1315	Attribute Name: expression_importance_score: 0.5
1316	candidates: [scared: 0.8, determined: 0.1, playful: 0.1]
1317	
1318	
1319	Dog
1320	Attribute Name: breed, importance_score: 0.8,
1321	Shepherd: 0.15, Bulldog: 0.1, Beagle: 0.1,]
1322	Attribute Name: coat_color, importance_score: 0.7,
1323	candidates: [brown: 0.2, black: 0.2, white: 0.2,]
1324	candidates: [long: 0.3, short: 0.3, fluffy: 0.2, shaggy: 0.1, wayy: 0.1]
1325	Attribute Name: hat, importance_score: 0.5,
1326	candidates: [baseball cap: 0.2, bowler hat: 0.2, top hat: 0.2,]
1327	
1328	
1329	Relations
1330	Dog-Rabbit
1331	importance score: 0.9,
1332	spatial_relation: [chasing: 1.0]
1333	Coat-Dog
1334	importance_score: 0.8,
1335	adornment_relation: [wearing: 1.0]
1336	
1337	



We employ the merged prompt across all agent variations (Ag1, Ag2, Ag3) to generate images at each turn of the interaction. Belief is being used in Ag1 and Ag2, which utilize belief parsing to extract structured representations from the user's input. Ag3 relies solely on the LLM's inherent ability to grasp the user's needs from the conversation history and merged prompt, without explicit belief state construction.

1356 D.2 IMPLEMENTATION OF ACTION

1355

1357

1363

1364

1365

1367

1369

1370

From an information theoretic perspective, an optimal action is the one that maximizes the information gain between the observation and the belief, i.e. $a_t = \arg \max_a H(o_{i-1}; b_{i-1} \mid a) - H(o_i; b_i \mid a)$. However, directly optimizing this objective can be computationally challenging. Therefore, we explore several heuristic strategies to effectively reduce uncertainty:

• Maximize the overall heuristic importance score (MHIS): This strategy focuses on maximizing the overall importance score of the entities, attributes, and relations within the belief. We further ask a question regarding an attribute or relation by maximizing the overall heuristic importance score. The score can be modeled as:

$$max_{e,a,c,r}(IS(\mathbf{e}) * IS(\mathbf{a}) * P(\mathbf{e}) * Ent(\mathbf{c}), IS(\mathbf{r}) * P(\mathbf{r}) * Ent(\mathbf{c}))$$
(1)

Here IS, P, Ent represents importance to ask score, probability of appearing, and entropy of the probabilities respectively and e, a, c, r represents entity, attribute, candidate list, relation respectively.

- Ask Important Clarification Question based on belief (*AICQ_B*): This strategy leverages the structured information within the belief. We provide the LLM with the user prompt, conversation history, and the current belief, utilizing an ICL prompt (§D.14) to guide question generation. The LLM then formulates a clarification question aimed at eliciting information about key features of the image, naturally prioritizing those with higher *Importance to ask score* within the belief.
- 1376
1377• Ask Important Clarification Question directly $(AICQ_{base})$: This strategy relies on the
LLM's inherent ability to identify important aspects of the user prompt and conversation
history. The LLM (§D.13) generates an important clarification question based on its implicit
understanding of the user's needs, without explicitly relying on the structured information
in the belief.

Ag1 employs MHIS strategy for question generation. This strategy leverages the importance scores assigned to entities, attributes, and relations within the belief state. It identifies the element with the highest heuristic importance score and formulates a question aimed at eliciting further information about that specific element. The question is then verbalized using the LLM described in Section §D.15.

Ag2 utilizes the parsed belief state as the basis for question generation. It employs the $AICQ_B$ strategy, which leverages the structured information within the belief state to generate targeted clarification questions.

 $\begin{array}{l} Ag3 \text{ relies solely on the conversation history for question generation. It employs the <math>AICQ_{base} \\ \text{strategy, which leverages the LLM's ability to understand the ongoing dialogue and identify key areas \\ \text{requiring further clarification.} \end{array}$

1393

1394 D.3 IMPLEMENTATION OF TRANSITION

1395 Both Ag1 and Ag2 require belief updating to incorporate new information gained Belief Updating. 1396 during the interaction. At each turn, we perform prompt merging to create a comprehensive prompt that summarizes the conversation history. This merged prompt is then used for belief parsing to 1398 obtain an updated belief state. For Ag2, this updated belief state directly informs the subsequent 1399 interaction. For AgI, it incorporates additional post-processing mechanisms to enhance memory and 1400 prevent redundant questioning: (i) Redundancy elimination: If an attribute or relation has already been addressed in the conversation history, the corresponding user response is assigned as the sole 1401 candidate with a probability of 1.0, and its importance score is set to 0. This prevents the agent from 1402 repeatedly asking about the same information. (ii) Information retention: If an attribute or relation 1403 from the conversation history is absent in the updated belief state, it is explicitly added. This ensures that the agent retains crucial information even if it's not explicitly present in the latest parsed belief
 state.

1408 D.4 USER SIMULATION

To simulate end-to-end agent-user interactions, we implement a user simulator that mimics human question-answering behavior. This simulator operates as follows:

- It generates a belief state based on a ground truth prompt, representing the user's intended image. This serves as the simulator's internal representation of the desired image.
- Mirroring the $AICQ_B$ strategy, the simulator takes the ground truth prompt, conversation history, and its current belief state as input. It then leverages an ICL prompt (see §D.14) to generate a response to the agent's question. This ensures that the simulator's answers are consistent with its internal belief state and the ongoing conversation.
- 1417 1418

1420

1407

1409

1410

1411 1412

1413

1414

1415

1416

1419 D.5 AG1: HEURISTIC SCORE AGENT

1421 The *Heuristic score agent* leverages the importance scores and probabilities within the belief state 1422 to guide its question-asking strategy. The underlying principle is to identify and inquire about the 1423 entity, attribute, or relation that exhibits both high importance and high uncertainty. This aligns with the uncertainty reduction principle discussed in §4.1.1, which emphasizes minimizing uncertainty 1424 through targeted questioning. To achieve this, we define a *heuristic importance score* as formulated 1425 in Equation 1, and the agent then selects the attribute or relation with the highest heuristic importance 1426 score as the focus of its inquiry. To facilitate easy answering, we utilize an LLM to generate user-1427 friendly questions with multiple-choice options. For example, the agent might ask: What color of 1428 the rabbit do you have in mind? a. black, b. white, c. brown. d. unkown. If none of these options, 1429 what color of the rabbit do you have in mind?. This format allows users to simply select the most 1430 appropriate option or provide their own answer if needed.

1431 Here's a summary of Ag1's implementation: (i) State Representation: The agent's state comprises 1432 the merged prompt and the current belief state. (ii) **Select Action**: *MHIS* strategy is employed to 1433 identify the attribute or relation of interest based on the heuristic importance score. (iii) Verbalize 1434 Action: An LLM (§D.15) is used to generate a clear and concise question about the selected attribute 1435 or relation. (iv) **Answer Ouestion**: The user simulator provides an answer to the agent's question, 1436 mimicking human response behavior. (v) Transition: The agent updates the merged prompt with 1437 the new information, re-generates the belief state based on the updated prompt, and applies the 1438 post-processing logic outlined in §D.3 to ensure consistency and prevent redundancy.

- 1439 1440
- 1441 D.6 AG2: BELIEF-PROMPTED AGENT

1442 The Ag2 agent incorporates the belief state into its decision-making process but adopts a different 1443 approach compared to Agl. Instead of relying on a heuristic score, Ag2 leverages the full capacity 1444 of an LLM to generate questions. It provides the LLM with comprehensive information, including 1445 the merged prompt, belief state, and conversation history, allowing the LLM to formulate the most 1446 informative questions possible. To guide the LLM towards generating effective questions, we 1447 incorporate specific instructions in the prompt, emphasizing the following principles: The question 1448 should be as concise and direct as possible. The question should aim to obtain the most information 1449 about the style, entities, attributes, spatial layout and other contents of the image. Remember to ask for information that are critical to knowing the critical details of the image that is important to the 1450 user. The question should reduce your uncertainty about the user intent as much as possible. 1451

Here's a summary of Ag2's implementation: (i) **State Representation**: The same as Ag1, the agent's state consists of the merged prompt and the current belief state. (ii) **Select Action**: $AICQ_B$ strategy is employed, which leverages an LLM to generate a question based on the comprehensive input information. (iii) **Verbalize Action**: Since the LLM directly generates the question, no separate verbalization step is required. (iv) **Answer Question**: The user simulator provides an answer to the agent's question. (v) **Transition**: The agent updates the merged prompt with the new information and re-generates the belief state based on the updated prompt.

1458 D.7 AG3: PRINCIPLE-PROMPTED AGENT

A simple and effective implementation of LLM-based multi-modal dialogue systems is to use the context to store the history of conversations between the system and the user, and directly generate the next response based on the context.

To align with the principles outlined in §4.1.1, we guide the LLM's question generation with a prompt (§D.13) that emphasizes all principles: *Based on the original prompt and chat history please provide a question to ask about the image. The question should be as concise and direct as possible. The question should aim to learn more about the attributes and contents of the image, the objects, the spatial layout, and the style.*. The prompt also includes the history of conversation. This strategy aims to generate questions that are easy for users to understand and answer, while effectively reducing the *agent's uncertainty about the desired image.*

1470Here's a summary of Ag3's implementation: (i) State Representation: The same as Ag1 and Ag2, the1471agent's state consists of the merged prompt and the current belief state. (ii) Select Action: $AICQ_{base}$ 1472strategy is employed, which leverages an LLM to generate a question based on the conversation1473history. (iii) Verbalize Action: The LLM directly generates the question, so no separate verbalization1474step is needed. (iv) Answer Question: The user simulator provides an answer to the agent's question.1475(v) Transition: The same as Ag2, the agent updates the merged prompt with the new information and1476re-generates the belief state based on the updated prompt.

1512 1513 D.8 ENTITY PARSER PROMPT INSTRUCTION

1514	1	Given a text-to-image prompt list out all the entities that are mentioned in the prompt.
1515	2 3	** Explicit Entities :** List all clearly stated entities within the prompt (people, objects, animals, locations, etc.).
1516	4	** Implicit Entities :** Identify potential entities that are implied or strongly suggested by the prompt, even if not explicitly mentioned.
1517	6	**Weather:** If the scene or mood suggests specific weather conditions (sunny, rainy, stormy, etc.).
1518	8	**Location:** If a general or specific setting is finited at (indoors, outdoors, a particular city/landscape, etc.). **Time of Day:** If the prompt implies a certain time (dawn, midday, dusk, night).
1519	9 10	**Mood or Atmosphere:** If the prompt evokes a particular emotion or ambiance (joyful, mysterious, peaceful, etc.).
1520	11	The estant should be list and estimate a ball be formered as a TGON list with the following fields
1521	12 13	The output should be list and each entry should be formated as a JSON dict with the following fields :
1522	14 15	"name": The name of the entity. "importance.to.ask_score": The importance score of asking a question about this entity to reduce the uncertainty of what the image is given the
1523		user prompt. Make sure that this is a number between 0 and 1, higher means more important. Consider these factors when assigning scores: 1.
1524		strongly influence the layout of the image, such as the position or portrayal of other entities in the scene; 3. significantly descrease the
1525		score for entities that are already well specified in the prompt; 4. significantly increase the score for implicit entities that are likely to appear in the image and their appearance can significantly impact the image.
1526	16 17	"description": A short description of the entity. "entity_type": The type of this entity. It could be either explicit, implicit, background. No other value is allowed.
1527	18	" probability of appearing ": The probability of the entity appearing in the image. This is a number between 0 and 1. You should assign a probability with the following rules in mind:
1528	19	1. If the prompt says an entity does not exist, assign a 0.0 probability. Because the entity does not exist, you should also assign 0 to
1529	20	importance_to_ask_score of this entity. 2. If the prompt indicates an entity definitly exists in the image, assign a 1.0 probability.
1530	21	3. If the prompt does not say anything about the existence of the entity, assign a probability between 0 and 1. This probability is higher if the entity is more likely to appear in the image given the context specified by the prompt
1531	22	4. If the prompt says an entity exists but there is an indication that the entity is not likely to appear in the image, assign a probability
1532	23	between 0 and 1, higher if the entity is more likely to appear in the image.
1533	24 25	Below is an example input and output pair: Example1:
1534	26	Input: {{
1535	28	user_prompt : generate an image of a nonneau rabot running on grass with sun simming. There is no uses in the background. }}
1536	29 30	Output: [{{
1537	31 32	"name": "rabbit ", " importance to ask score ": 0.5
1538	33	"description": "a lionhead rabbit",
1539	34 35	" entity_type ": " explicit ", " probability_of_appearing ": 1.0
1540	36 37	}}, {{
1541	38	"imme": "grass ",
1542	39 40	"description ": "grass",
1543	41 42	"entity_type": "explicit", " probability_of_appearing ": 1.0
1544	43 44	}}, {{}}
1545	45	"name": "sun",
1546	40 47	"description ": "sun is shining",
1547	48 49	" entity_type ": " explicit ", " probability_of_appearing ": 0.3
1548	50 51	}}, //
1549	52	"name": "sun light ",
1550	53 54	" importance.to_ask_score ": 0.1, " description ": "sun light shining on the grass and the rabbit ",
1551	55 56	" entity_type ": " explicit ", " probability_of_appearing ": 1.0
1552	57	}}, {{
1553	58 59	11 "name": "tree ",
1554	60 61	" importance.to_ask_score ": 0, " description ": " trees in the background",
1555	62 63	" entity_type ": " explicit ",
1556	64	<pre>}} }</pre>
1557	65 66	۱۱ "name": "camera angle",
1558	67 68	" importance.to.ask_score ": 0.8, " description ": "the camera angle of the image",
1559	69 70	" entity_type ": "background",
1560	70)}},
1561	72 73	{{ "name": "weather",
1562	74 75	" importance_to_ask_score ": 0.8, " description ": "weather".
1563	76	" entity_type ": "background",
1564	77 78	probability_of_appearing : 1.0 }},
1565	79 80	{{ "name": "image style ".
	81	" importance.to.ask.score ": 1.0,

1566		
1567	82 83	" description ": "the style of the image", " entity type ": "background"
1568	84	" probability_of_appearing ": 1.0
1560	85 86	}}, {{
1570	87 88	"name": "background color", "importance to ask score": 0.8
1570	89	"description": "the background color of the image",
1571	90 91	" entity_type ": "background", " probability_of_appearing ": 0.5
1572	92 03	}}
1573	93 94	J
1574	95 96	[[a few additional examples]]
1575	97	Identify the entities given the input given below. Strictly, stick to the format
1576	99 99	Input: {{
1577	100 101	"user_prompt": "{user_prompt}" }}
1578	102	Output:
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		
1612		
161/		
1615		
1610		
1010		
101/		
1018		
1019		

1620 D.9 Attribute Parser Prompt Instruction

_

1622	1	Given a text-to-image prompt and a particular entity described in the prompt, and your goal is to identify a list possible attributes that could describe the particular entity. Output Requirements:
1623	2	econo de partenar entry. Oupur requirementais.
1624	3 4	 If this attribute has already existed as an entity in other existing entity list, then do not include it. the attribute candidate could be a mixed of values like 'color A and color B'.
1625	5	3. The output should be a json parse-able format:
1626	7	name (str): The name of the attribute.
1627	8	importance_to_ask_score (float): The importance score of asking a question about this attribute to reduce the uncertainty of what the image is given the user prompt. This is a number between 0 and 1, higher means more important. Consider these factors when assigning scores: 1.
1628		Increate the score for attributes that are the primary attributes of an important entity; 2. significantly increase the score for attributes that could strongly influence the generation or portrayal of OTHER attributes in the scene; 3. descrease the score for
1629		attributes that are already well specified in the prompt. For example, a breed of a dog would impact other attributes like color, size, etc
1630		. So the breed attribute should have a higher importance score than color, size, etc. Assign a much lower score if the attribute's value is already mentioned in the user prompt.
1631	9	candidates (List of names and probabilities): List of possible values that the attribute can take. Make sure to generate atleast 5 or more possible values. These should be realistic for the given entity. For each attribute, returns the probability that the user wants this candidate
1632		based on the user prompt. If it's already mentioned by the user, only generate one candidate (the mentioned one) and assign 1.0 as the
1633		dog with breed Samoyed, the color attribute has a very high probability of white.
1634	10 11	Below are two examples of input and output pairs :
1635	12	Example 1
1636	14	Input: {{
1637	15 16	user_prompt : generate an image of a white rabbit running on grass , "entity ": "rabbit ",
1638	17 18	" other_existing_entities ": "grass"
1639	19	Output: [
1640	20	11 "name": "color",
1641	22 23	" importance.to_ask_score ": 0.9, "candidates ": {{"white":1.0}}
1642	24 25	}},
1643	26	"name": "breed",
1644	27	"candidates ": {{"Dutch": 0.20,
1645	29 30	"Mini Lop": 0.15, "Netherland Dwarf": 0.15,
1646	31	"Lionhead": 0.10, "Element Giane": 0.10
1647	33	"Mini Rex": 0.10,
1648	34 35	"English Angora": 0.08, "Mini Satin": 0.05,
1649	36 37	"Himalayan": 0.05, "Californian": 0.02}}
1650	38	}}, }},
1651	39 40	11 "name": "age",
1652	41 42	" importance.to.ask_score ": 0.1, " candidates ": {{"adult": 0.6,
1653	43 44	"baby": 0.2, "senior": 0.233
1654	45	<pre>}}</pre>
1655	46 47	
1656	48 49	[[a few additional examples]]
1657	50 51	Generate attributes given the input given below. Do not include other entities in the attributes. Strictly stick to the format.
1658	52	"user_prompt": "{user_prompt}",
1659	53 54	" other_existing_entities ": "{ existing_entities }"
1660	55 56	}} Output:
1661		•
1662		
1663		
1664		
1665		
1666		
1667		
1668		
1669		
1670		
1671		
1672		
1673		

1674 D.10 RELATION PARSE PROMPT INSTRUCTION

10/5		
1676	1	Given a text-to-image prompt and a list of entity described in the prompt, your goal is to identify a list of entity pairs that have relations between them. Ignore entity pairs without relations. The output should be a json parse-able format (No comma after the last element of the
1678	2	nst):
1679	3	Input: user_prompt: the prompt from the user.
1680	5 6	entities : a list of entities mentioned in the user_prompt.
1681	7	Output: name (str): The name of the relation. Use 'entity1-entity2' as the format.
1682	9 10	description (str): A short description of the relation . spatial relation (man from potential relation candidates to probability). Possible spatial relations between the two entities. If a relation is
1683	11	mentioned in the user prompt, assign 1.0 as the probability. The sum of probabilities over all relation candidates shall be 1.
1684	11	and 1, higher means more important. Assign a higher score if the two entities are very important, the relation between them is very unclear and the relation is very important for the layout of the image
1685	12	name.entity.1 (str): The name of the first entity.
1686	13 14	name_entity_2 (str): The name of the second entity . is_bidirectional (bool): Whether the relation is bidirectional .
1687	15 16	Below is an example input and output pair:
1688	17	Example1:
1689	19	"user-prompt": "generate an image of a lionhead rabbit sitting on grass, and a eagle is flying through the sky. There is a tree in the
1690	20	background.", "entity ": [" rabbit ", "grass", "eagle", "tree "]
1691	21 22	}} Output: [
1692	23 24	{{ "nama": "rabhit_arree"
1693	25	"description": "rabbit grass",
1694	26 27	" spatial_relation ": {{ "above": 0.8, "below": 0.0, "in front of": 0.0, "behind": 0.0, "left of": 0.1, "right of": 0.1}}, "importance_to_ask_score ": 0.1,
1695	28 29	"name.entity_1":" rabbit ", "name.entity_2": "grass".
1696	30	" is_bidirectional ": true
1697	32	
1698	33 34	"name": "eagle-grass", "description ": "eagle is flying through the sky",
1700	35 36	" spatial_relation ": {{"above": 1.0, "below": 0.0, "in front of": 0.0, "behind": 0.0," left of ": 0.0, "right of": 0.0}}, "importance_to_ask_score ": 0.1,
1700	37	"name_entity_1":"eagle",
1701	39	is_bidirectional ": false
1702	40 41	}}, {{
1704	42 43	"name": " tree – grass", " description ": "",
1705	44 45	" spatial_relation ": {{"above": 0.5, "below": 0.0, "in front of": 0.0, "behind": 0.0, "left of": 0.25, "right of": 0.25}}, "importance to ask score ": 0.1
1706	46	"name.entity.1.": "tree",
1707	48	is_bidirectional ": false
1708	49 50	}},
1709	51 52	[[a few additional examples]]
1710	53 54]
1711	55	Identify relationships between entities given the input given below. Strictly stick to the format.
1712	56 57	Input: {{ "user_prompt": "{user_prompt}",
1713	58 59	"entity ": "{entity_names}" }}
1714	60	Öutput:
1715		

1716 D.11 VERBALIZATION PROMPT INSTRUCTION

1 The chat history is as follows:

1717 1718 1719

3

- 1720
- 1721 1722
- D.12 MERGE PROMPT PROMPT INSTRUCTION

1702	
1723	1 You are writing a prompt for a taxt-to-image model based on user feedback. The original prompt is {prompt}. The user has provided some additional
4 7 0 4	Tot are writing a prompt for a text-to-image model based on user reedback. The original prompt is {prompt is {prompt j. The user has provided some additional
1724	information: { additional_info }. Please write a new prompt for the text-to-image model. The new prompt should be a meaningful sentence or a
1705	paragraph that combines the original prompt and the additional information. Do not add any new information that is not mentioned in the
1720	prompt or the additional information. Make sure the information in the original prompt is not changed. Make sure the additional information
1726	is included in the new prompt. Make sure the new prompt is a description of an image. If the additional information or the original prompt
1120	specifically says that a thing does not exist in the image, you should make sure the new prompt mentions that this thing does not exist in
1727	the image. DO NOT generate rationale or anything that is not part of a description of the image.

question : {action. verbalized.action } and answer: {observation}. Turn the question and action into a single declarative sentence that describes the answer – do not phrase it as a question. Example output: the firetruck in the image is red.

1730	1	[[Instruction for the first question]]					
1731	2	The original prompt was: {self.original.prompt } - Based on the original prompt please provide a question to ask about the image. The question					
1732		should be as concise and direct as possible. The question should aim to learn more about the attributes and contents of the image, the objects, the spatial layout, and the style. Make sure that you question the answer within <question> and </question> markers					
1733	4	[[Instruction for the following question]]					
1734	6	saved on the chat history please provide a new question to ask about the image the chat history is as follows and is enclosed in <i>chat history</i>					
1735	/	Based on the chat history please provide a new question to ask about the image, the chat history is as follows and is enclosed in <chat.history> and </chat.history> markers:{self. chat.history} The question should be as concise and direct as possible. The question					
1736		should aim to learn more about the attributes and contents of the image, the objects, the spatial layout, and the style. Make sure that you question the answer within <question> and </question> markers.'					
1737	ı						
1738		D 14 $AICO_{\rm T}$ Prompt Instruction					
1739		D.14 AICQBINOMETINSINCETION					
1740	1	You are an intelligent agent that helps users generate images. Before generating the image requested by the user, you should ask the most important					
1741	2	The user describes the image as: {user_prompt}.					
1742	3	The following is your belief of what the image contains, including the entities, attributes of each entity and relations between entities. Each entity has "name", "descriptions", "importance to ask score" and "probability of appearing". "Name" is the identifier of the entity."					
1743		Descriptions" is the description of the entity. "Importance to ask score" is how important it is for the agent to ask whether the entity exists. Probability of appearing" is the probability the agent estimated that this entity exits in the image.					
1744	5	Fash antity has a list of attributes. Each attribute has "none" "immortance to adv energy" and "condidates" "None" in the identifier of the					
1740	0	attribute . "Importance to ask score" is how important it is to ask about the exact value for the attribute of the entity. "Candidates" is a					
1740	7	list of possible values for the attribute.					
17/19	8 9	Each candidate value has a probability that describes how likely this candidate value should be assigned to the attribute . For example, "Attribute Name: color, Importance to ask Score: 0.9, Candidates: [white: 0.5, black: 0.5]" means the color is either white or black.					
17/10	-	a chample, refine rather contrained to as door, or, candidades, while or, other or, internet the contrained with of black, each with 0.5 probability. If you ask about attributes, you should ask about the attributes with the highest uncertainty. Your uncertainty					
1750	10	fairly certain.					
1751	10 11	The agent belief is:					
1752	12 13	{ belief_statestr () }					
1753	14	Based on the user prompt "{user_prompt}" and the belief of the agent, please provide a question to ask about the image. The question should be as concise and direct as possible. The question should aim to obtain the most information about the style entities attributes spatial					
1754		layout and other contents of the member to ask for information that are critical to knowing the critical details of the image that					
1755		is important to the user. The question should reduce your uncertainty about the user intent as much as possible. DO NOT ask question that can be answered by common sense. DO NOT ask question that are obvious to answer based on the user prompt "{user_prompt}". DO NOT ask any					
1756	15	question about information present in the following user-agent dialogue within <dialogue> and </dialogue> markers.					
1757	16 17	<dialogue> {conversation}</dialogue>					
1758	18						
1759	20	DO NOT ask any question that has been asked in the dialogue above.					
1760	21 22	Your question does not have to be entirely decided by the belief . You can construct any question that make yourself more confident about what the					
1761	23	image is. Think step by step and reason about your uncertainty of the image to generate. Make sure to ask only one question. Make sure it is not very					
1762		difficult for the user to answer. For example, do not ask a very very long question, which can take the user a long time to read and answer					
1763	24	Make sure that you question the answer within <question> and </question> markers.					
1764							
1765		D.15 HSA OUESTION PROMPT INSTRUCTION					
1766	r						
1767	1 2	You are constructing a text-to-image (T2I) prompt and want more details from the user. You have to ask a question about the the most important entity or the attribute of the most important entity.					
1768	3	We have entity types: (i) explicit: directly ask question with options; (ii) implicit: ask whether this entity required for the image with yes or no as options; (iii) background; impore the attribute value and directly ask the value of the entity. (iv) relation; add keyword like '					
1769	4	relation 'to emphasize this entity is a relation.					
1770	4	question and follow it with options.					
1770	5 6						
1772	7 8	Example1: entity: rabbit					
177/	9 10	attribute : color					
1775	11	entity_type : explicit					
1776	12	question: what color of the rabbit do you have in mind? a. black, b. white, c. brown. d. unkown. If none of these options, what color of the rabbit do you have in mind?					
1777	13 14	[[a few additional examples]]					
1778	15 16	Example5:					
1779	17	entity: Sentity					
1780	18 19	candidates : \$candidates					
1781	20 21	entity_type : \$entity_type question :					
	l						

1.	Input: Initial prompt n. User a Agent a (with n.) man turna
1.	input . Initial prohibit p_0 , Osel <i>a</i> , Agent <i>a</i> (with p_0), <i>max_tarns</i>
2:	Output: Refined prompt p_f
3:	$p_f \leftarrow p_0$
4:	for $turn_i d = 0$ to $max_t urns - 1$ do
5:	$action \leftarrow a.SelectAction()$
6:	$question \leftarrow a. VerbalizeAction(action)$
7:	$answer \leftarrow u.AnswerQuestion(question)$
8:	a.Transition(action, answer)
9:	$p_f \leftarrow a.prompt$
10:	end for
11:	return p_f

1797

1796 E BACKGROUND

T2I Evaluation. Inception Score (Salimans et al., 2016) and Frechet Inception Distance (Heusel et al., 2018) are popular metrics to measure the fidelity of generated images, i.e. the similarity of 1799 the generated images to real ones. Improved precision and recall (Kynkäänniemi et al., 2019) allows to analyse the sample quality and the coverage independently. Since text prompts are used to guide 1801 image generation in T2I models, image-prompt alignment is an important evaluation metric which can 1802 be classified as embedding-based such as CLIPScore (Hessel et al., 2022), ALIGNScore (Zha et al., 1803 2023), VQA-based metrics such as TIFA (Hu et al., 2023), DSG (Cho et al., 2023) abd VQAScore (Lin et al., 2024) and captioning based metrics like LLMScore (Lu et al., 2023). Approaches such 1805 as PickScore (Kirstain et al., 2023), ImageReward (Xu et al., 2023) and HPS-v2 (Wu et al., 2023) finetune models on human ratings to devise a metric that aligns with human preferences. Recently, 1807 diversity of generated images (Naeem et al., 2020) is becoming an important metric of measurement to track progress, especially in the geo-cultural context (Kannen et al., 2024; Hall et al., 2024). 1808 Diversity can be used to quantify the under-specification in the input prompt: more specific the 1809 prompt, the less diverse are the generated set of images across different seeds (Kannen et al., 2024). 1810

Prompt expansion is a widely known technique to improve image generation (Betker et al., 2023).
ImageinWords (Garg et al., 2024) proposes to obtain high-quality hyper-detailed captions for images, which significantly improve quality of image generation. Datta et al. (2024) present a generic prompt expansion framework used along Text-to-Image generation and show an increase in user satisfaction through human study. Our work can be viewed as a method to adaptively expand a T2I prompt based on user feedback. Samples from the agent belief can be used to construct expanded prompts.

1817 1818

F AUTOMATED EVALUATION

1819 1820

In Algorithm 2, we show the user-agent self-play procedures that we used to perform all automated evaluation.

1821 1822

1824

G DETAILS ON THE AGENT INTERFACE

Below is a showcase of how users could interact with the belief graph and clarifications in a hypothesised interface, to better iterate their inputs, to reach higher a quality and satisfaction of outputs. This is a crudely hypothesised, intentionally simple interface for the sake of research, but could be iterated and improved upon in many ways depending on application and users.

1. Default state On load of the app, there would be a text prompt input and space for output images, as is common across typical T2I interfaces. There would also be space for the user to view either clarifications from the model, or a graph interface, as part of the overall "input" section as these would act as a further input for future model output iterations. See Figure 11 below as reference.

2. Output images, with Clarifications Once the user has submitted the prompt and the model has responded, there would be a set of images, as initial outputs from the users prompt. Below the input prompt would be a set of "Clarifications" in its populated state. These clarifications would ask the



Figure 11: Default state of a possible interface.

user specific questions that would be necessary to increase the specificity of the prompt, for the model to get a more accurate results aligned to the users intention, or to help the user realise their intention.

Options would be given of the highest probability options for each Clarification, but the user could also fill in a totally new option via a free text field. Once answered by selection or text input, the clarifications would be added to the above, primary prompt for regeneration when the user selects. See Figure 12 below as reference.



Figure 12: Interface once prompt has been input with clarifications.

3. Graph Entities & Attributes Instead of the clarifications, the user could select to instead view 1917 a Graph by clicking Tab above the clarifications themselves. This graph would be populated will 1918 all Entities from the prompt explicit and implicit visually defined differently (in this diagram by 1919 the dotted line surrounds implicit entities, but is a filled line when surrounding explicit entities). 1920 The graph layout will be structured, depicting relationships concentrically i.e. "on", "in" or "under" 1921 for example, will become child entities, and be displayed within the parent entities' boundary. For 1922 example a 'Mug' that has the relationship of 'on' a 'Table' entity, will sit within the boundary of 1923 'Table', as also would a 'Plate' if that had the same child-parent relationship. 1924

Below the Graph would also be a list of 'cards' (i.e. boxed groups of information), one for each "explicit" or "implicit" entity. Within each card a user could see the status of implicit / explicit, and change this status to confirm or deny its presence. The user could also see a list of "attributes" associated to that entity, which the model has assumed. Each of these attributes could be changed by interacting with a list of alternatives via drop down. These lists are determined in terms of which items and order of items, based on the probability by which the model sees them, ordered with higest first. This probability would be made clear to the user to define the order by seeing the peercentage next to the label. See Figure 13 below as reference.

4. Graph Relationships The user would also be able to change the state of the Graph and Cards, to instead focus on the relationships between entities, by toggling to "Relations". In this state the user would be able to focus on two specific entities (e.g. 'mug' and 'table'), see the description of the relationship (e.g. 'the mug is sitting on the table') and if desired change the relationship to an alternative (e.g. 'on', changed to 'under') via a drop down of options which the model determined as alternative options ordered by probability, as per attributes. See Figure Figure 14 below as reference.

Once any of these changes are made the user could initiate a regeneration via the updated prompt to create a new set of output images, which can then be further refined via the same method.

1941

1915

1916

1942 H DETAILS ON USER STUDY

1943

Below we describe the exact guideline definitions we shared with the user for a user study.

INPUT		DUTPUTS	
A dominating image of the Eiffel prominently displayed on the str	Tower with the Olympic rings ucture.		
CLARIFICATIONS GRAPH			
Myseaf Interes	Backgrund extitles Disputine Stay R Cloudy U U U U U U U U U U U U U	Image goes here	Image goes he
View: Control 5 a stationars V	Seine Norm Seine Norm eprev. 1 of 3 <u>necto</u> Athlocate Markowski Million Markowski Million Markowski Million Norm Million Markowski Million Norm Million Markowski Million	Image goes here	Image goes fe

Figure 13: Interface with Graph displaying Entities, with cards below enabling a user to change attributes associated to each entity.

🛱 Арр			
INPUT		OUTPUTS	
A dominating image of the Eiffel Tower prominently displayed on the structure	with the Olympic rings		
CLARIFICATIONS GRAPH			
Physical Entities	Background entities	Image goes here	Image goes h
'Paris Cityscape'	'Daytime'		
Tower' 1 Sky'	'Camera angle'		
('Olympic Rings' 'Clouds' Strakoti	'Lighting'		
	/		
12 People' 14 Trees' 15 'Sei	ne River'		
View: () Entities & attributes ↔ Relationships	<prev: 1="" 2="" <u="" of="">next></prev:>		
Relation: Eliffel Tower - Olympic Rings Description: "Olympic rings prominently displayed near the top of the Eliffel Tower" Span	ation: Eliffel Tower - Paris Cityscape sripsor: "The Eliffel Tower in the city of Paris" tal relationship: ("On Champs da Mars" -		
		Image goes here	Image goes he
(n.) Heristoon Parts Untyscoper * Sky Description "The ky above Paris" Sparial relationship: Clearly above" *	scon: Sky * Leouas cription: "Clouds in the sky above Paris" fai relationship: "Spread evenly across" *		
(12) Relation: Paris Cityscape - People	ation: Paris Cityscape - Trees		
Description: "The papepte visible anound Paris streets" Spatial relationship: "Orowsded streets" • Spatial relationship:	cription" the trees visible around Paris scenery" fail relationship: "Sparsoly scattered" •		

Figure 14: Interface with Graph displaying relations between Entities, with cards below enabling a user to change relationships between entities.

H.0.1 HYPOTHESIZED FRUSTRATIONS

We presented participants with the following hypothesized frustrations related to T2I model usage:

- 1. Prompt Misinterpretation: The model misunderstands complex relationships between entities in the input prompt.
- 2. Many Prompt Iterations: The model does not immediately generate what the user intends, requiring numerous iterative changes to the input prompt.

1998 1999 2000	3. Inconsistent Generations: The model reinterprets the input prompt differently between iterations, causing unwanted changes in the generated images.
2001 2002 2003 2004	4. Incorrect Assumptions: The model makes incorrect assumptions or no assumptions when encountering gaps in the details provided in the input prompt, leading to undesired outputs.
2005	Explanations of terms were given to users of:
2007 2008 2009	1. "Entities" are single items that are intended to be in the image e.g. "Cat" and "Ball", from "make a sketch of a Cat playing with a Ball"
2010 2011 2012	2. "Prompt" means the text written to communicate the intended output image e.g. the sentence "make a sketch of a Cat playing with a Ball" is the "Prompt", also known as "Input"
2013 2014 2015	3. "Iterations" are each set of different image outputs by the model, taken from a different input, or even the same input just regenerated
2016 2017 2018 2019 2020	The question asked for each Frustration were: "Please score the below frustrations (or issues) that could be related to Text to Image AI Generation"."Rank in terms of how much they relate to your current usage, with your most commonly used model or app."
2021 2022	H.0.2 Hypothesized Features
2023 2024 2025	We proposed the following features as potential solutions to address the identified frustrations:
2026 2027 2028 2029	1. Clarifications: The model would ask specific clarifying questions about uncertainties in the prompt. These details would then be incorporated into subsequent iterations. For example: "Is the cat playing with: 1. a ball of wool, or 2. a tennis ball?"
2030 2031 2032 2033 2034	2. Graph of Prompt Entities: A visual representation of all entities in the prompt as a graph, allowing users to see and edit attributes of each entity. E.g., seeing that the model has assigned "round," "small," and "wooden" as attributes to "table" and allowing the user to change them to "square" and "metal."
2035 2036 2037 2038	3. Graph of Prompt Relationships: A visual representation of relationships between entities in the prompt, allowing users to see and edit these relationships. E.g., seeing that "donut" is "next to" "coffee" and allowing the user to change the relationship to "on top of."
2039 2040	The questions asked for each feature were:
2041 2042 2043 2044 2045	 "How likely this feature is to help your current workflow if you had it now?". With response options of: "Very unlikely to help", "Unlikely to help", "Could help", "Likely to help", "Very likely to help".
2046 2047	 "How soon would this feature deliver value to your work?" with response options of: "Very soon / immediately", "Sometime, "Not very soon".
2048 2049 2050 2051	Image references were given for each Feature as listed out below:

1. Clarifications:



Figure 15: Stimulus image in the survey to test the Model clarifications feature.

2. Graph of Prompt Entities:



Figure 16: Stimulus image in the survey to test the Model Graph of Entities and Attributes feature.

3. Graph of Prompt Relationships:



Figure 17: Stimulus image in the survey to test the Model Graph of Entity Relations feature.

H.0.3 HUMAN STUDY RESULTS

Table 3: Breakdown of the T2I usage frequency of the 143 participants recorded

Usage Frequency	No. of users	(%)
Many times a day	13	9.1
Many times a week	44	30.8
At least once a week	36	25.2
At least once a month	50	35.0

Table 4: Reported User Frustrations with existing T2I processes (% of users)

Frustration	V. Freq. (%)	Freq. (%)	Occas. (%)	V. Occas. (%)	No Issue (%)
Prompt Misinterpret.	7	19.6	43.4	23.1	7
Many Iterations	10.5	44.8	28	11.9	4.9
Inconsistent Gen.	11.2	20.3	39.9	21	7.7
Incorrect Assumptions	7	23.1	39.2	20.3	10.5

Table 5: Expected speed of value delivered from features (% of users)

Feature	Very soon / immediately (%)	Sometime(%)	Not very soon. (%)
Clarifications	57.7	37.2	5.1
Entity Graph	49.6	34.8	15.6
Relation Graph	41.8	44	14.2