# Multimodal Procedural Planning via Dual Text-Image Prompting

**Yujie Lu[1], Pan Lu[2], Zhiyu Chen[1], Wanrong Zhu[1], Xin Eric Wang[3]**
**William Yang Wang[1]**

[1]University of California, Santa Barbara, CA, USA
{yujielu,zhiyuchen,wanrongzhu,wangwilliamyang}@ucsb.edu
[2]University of California, Los Angeles, CA, USA
lupantech@gmail.com
[3]University of California, Santa Cruz, CA, USA
xwang366@ucsc.edu

**Task:** How to make traditional szechuan chicken?



**Step 1:** Gather the ingredients such as chicken, garlic, ginger, and light soy sauce.

**Step 2:** Marinade the chicken with a teaspoon of light soy sauce and cooking wine.

**Step 3:** Heat up a wok over medium-high heat until it starts to smoke and add oil.

**Step 4:** Add in marinated chicken pieces and stir fry until they turn golden brown.

**Step 5:** Add other ingredients like sichuan peppercorn powder and chili bean paste, then stir fry.

**Step 6:** Pour in remaining marinade mixture into the same wok and stir fry until combined.

**Step 7:** Return chicken pieces into the wok and stir fry until all liquid forms a glossy glaze.

**Step 8:** Transfer Szechuan Chicken onto serving plate garnished with green & serve immediately!
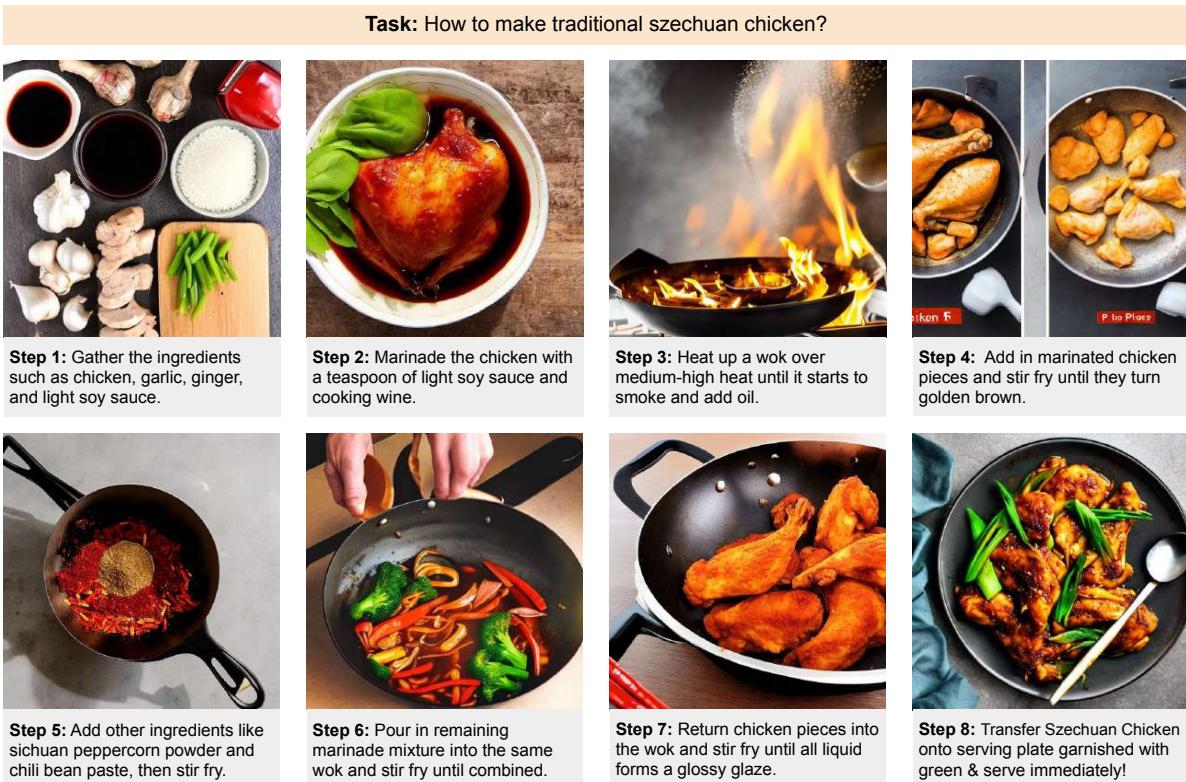
Figure 1: Our dual Text-Image Prompting (TIP) model generates coherent and authentic *multimodal procedural plans* with multiple steps towards a high-level goal, providing useful guidelines in task completion.

## Abstract

Embodied agents have achieved prominent performance in following human instructions to complete tasks. However, the potential of providing instructions informed by texts and images to assist humans in completing tasks remains underexplored. To uncover this capability, we present the multimodal procedural planning (MPP) task, in which models are given a high-level goal and generate plans of paired text-image steps, providing more complementary and informative guidance than unimodal plans. The key challenges of MPP are to ensure the informativeness, temporal coherence, and accuracy of plans across modalities. To tackle this, we propose Text-Image Prompting (TIP), a dual-modality prompting method that jointly leverages zero-shot reasoning ability in large language models (LLMs) and compelling text-to-image generation ability from diffusion-based models. TIP improves the interaction in the dual modalities using Text-to-Image Bridge and Image-to-Text Bridge, allowing LLMs to guide the textual-grounded image plan generation and leveraging the descriptions of image plans to ground the textual plan reversely. To address the lack of relevant datasets, we collect WIKIPLAN and RECIPEPLAN as a testbed for MPP. Our results show compelling human preferences and automatic scores against unimodal and multimodal baselines on WIKIPLAN and RECIPEPLAN in terms of informativeness,

arXiv:2305.01795v1 [cs.CL] 2 May 2023

temporal coherence, and plan accuracy.[1]

# 1 Introduction

Recent advances in embodied (Huang et al., 2022; Anderson et al., 2018) and conversational (Qiu et al., 2021) agents achieve prominent performance in task completion as humans by following instructions informed by texts and images. However, to what extent the models can provide useful guidelines for humans to complete the task remains underexplored. To uncover this, we propose the multimodal procedural planning task (as shown Figure 1). The task aims to generate goal-conditioned (e.g. "How to make traditional szechuan chicken") text (e.g. "a teaspoon of light soy sauce" explain how to marinade chicken in Step 2) and image (e.g. help identify the ingredients "chicken, garlic, ginger, and light soy sauce" in Step 1) plans as useful guidelines to assist humans in task completion.

Previous work (Huang et al., 2022) has explored the generation of procedural plans in text-only form. In contrast, we generate both text and image plans, which provide guidance for the agent to perform tasks that acquire complementary information from multimodal contexts. Generating plans in both text and image form poses new challenges since the generated plans should: a) be *informative* enough in both the text and image modalities, b) obey commonsense temporal *coherence*, such as the order of steps, and c) achieve high plan *accuracy*, indicating the complementary and alignment among multimodal plans.

Despite significant progress (Kojima et al., 2022; Song et al., 2022) in the development of large language models (LLMs), they are unable to generate images. Existing text-to-image (T2I) models can generate high-quality images conditioned on textual instructions (Ramesh et al., 2022; Rombach et al., 2022; Brooks et al., 2022). However, they are limited in their ability to generate images that require complex text comprehension, such as temporal reasoning (e.g. "learn basic surf safety *before* hitting the waves") and physical reasoning (e.g. "*pick up* the wine glass"). Additionally, generating text and image plans separately using LLMs and T2I models results in inconsistency and incoherence between the two modalities.

In this paper, we propose Text-Image Prompting (TIP), a novel dual-modality prompting framework that jointly leverages the capabilities of LLMs and

T2I models for multimodal procedural planning. We first generate vanilla text plans by directly asking LLMs (Kojima et al., 2022) for step-by-step procedures. To generate textual-grounded image plans, we devise the Text-to-Image Bridge (T2I-B), which elicits the complex language comprehension abilities of LLMs to assist T2I models in generating informative image plans conditioned on text plans. Similarly, we generate visual-grounded text plans using the Image-to-Text Bridge (T2I-B), which verbalizes the image plans and injects them back into LLMs to aid in revising the text plans, thereby improving their informativeness. The temporal coherence of the generated plans is also improved considering the context of both text and image. Benefiting from our dual-modality prompting, our generated plans are complementary and aligned across text and image modalities.

To address the lack of suitable datasets for evaluating multimodal procedural planning, we collect the WIKIPLAN and RECIPEPLAN datasets for benchmarking the task. We empirically evaluate the effectiveness of TIP on WIKIPLAN and RECIPEPLAN in a zero-shot setting and compare it with various baselines. Our results demonstrate that TIP generate plausible multimodal plans that are informative, temporally coherent, and accurate. Our work highlights the potential of combining knowledge from LLMs and T2I models to uncover multimodal zero-shot planning capabilities. Our main contributions are as follows:

- We introduce the multimodal procedural planning task and evaluate model performance using our collected WIKIPLAN and RECIPEPLAN datasets.
- We propose Text-Image Prompting (TIP), a dual-modality prompting approach that elicits procedural knowledge jointly from LLMs and T2I models, enabling visual-grounded text plans and textual-grounded image plans.
- We show that TIP substantially improves performance in terms of textual and visual informativeness, temporal coherence, and plan accuracy on human and automatic evaluations.

# 2 Related work

**Procedural Planning** This task (Zhang et al., 2020; Chang et al., 2020) has gain much attention in various aspects, including robotics (Tellex et al., 2011; Jansen, 2020; Brohan et al., 2022), vision-and-language navigation (Anderson et al.,
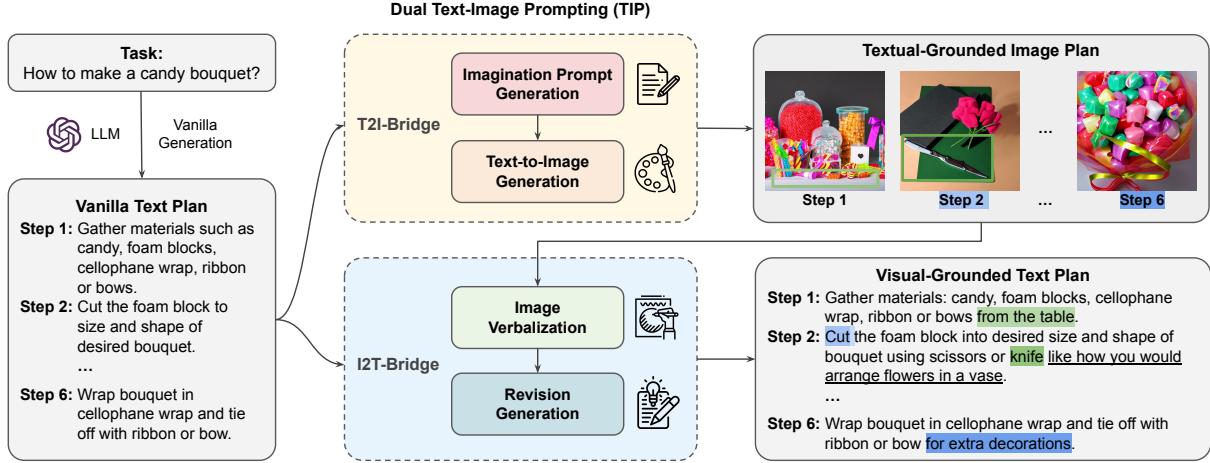
Figure 2: The vanilla text plan is generated using LLM. Our Text-Image Prompting (TIP) generates the textual-grounded image plan using T2I-Bridge (Fig. 3) and the visual-grounded text plan using I2T-Bridge (Fig. 5). The colors blue and green highlight the improved grounding in text and image respectively.

2018), conversational assistants (Ilievski et al., 2018; Qiu et al., 2021, 2022; Yang et al., 2022a), and animation (Zhao et al., 2022). Recent work is extended to the multimodal scenarios (Wu et al., 2022; Song et al., 2022; Wang et al., 2022c). In this work, we explore the multimodal procedural planning that generates goal-conditioned text and image sequences grounded in a multimodal context.

**Multimodal Generative Models**  Recently advanced diffusion models (Ramesh et al., 2022; Rombach et al., 2022) have shown remarkable abilities in generating high-quality images given text prompts. However, generating images with desired semantics requires proper prompts, which often come from a number of trials and errors (Liu and Chilton, 2022). To get more controllable generations, researchers have used large language models (LLMs) to expand input prompts with rich contextual knowledge. InstructPix2Pix (Brooks et al., 2022) combines the knowledge of GPT-3 and Stable Diffusion to generate large-scale examples of image editing as training data. In turn, recent advances in large-scale models based on transformers (Li et al., 2022; Wang et al., 2022a) exhibit incredible ability in image captioning, describing the given image using natural language.

**Injecting Visual Knowledge in LLMs**  Incorporating visual knowledge into large language models through visual imagination is a promising area of research. This can be achieved through the use of existing images as augmented visual features for language models, or through the generation of images to provide additional visual supervision to

language models (Yang et al., 2022b). Studies such as (Zhang et al., 2021b; Yang et al., 2022b; Zhu et al., 2022a; Lu et al., 2022b; Liu et al., 2022) have demonstrated the effectiveness of this approach. Our proposed TIP exploits the image descriptions in language form to inject the visual knowledge into LLMs and elicit its potential zero-shot reasoning ability to ground the textual sentences in the verbalized visual context.

## 3  Our Approach

### 3.1  Problem Definition

We formulate multimodal procedural planning as a conditional text and image sequence generation problem. Given a high-level goal $\mathcal{G}$ in natural language form, the model generates a sequence of low-level steps $\mathcal{S} = \{s_1, s_2, ..., s_n\}$. Each step $s_i$ in the sequence is represented by a paired text $t_i$ and image $v_i$ at timestep $i$. The text plan $\{t_1, t_2, ..., t_n\}$ and image plan $\{v_1, v_2, ..., v_n\}$ are both intended to be informative in their respective modalities and complementary across modalities. The final multimodal procedural plans ($\mathcal{S}$) is the combination of the text plan and image plans, which describe the procedure of completing the high-level goal.

### 3.2  Method Overview

We first elicit the zero-shot step-by-step reasoning ability in large language models (LLMs) to generate a vanilla text-only plan (left part in Figure 2). To enable grounding in multimodal context, we propose Text-Image Prompting (TIP), a dual-modality prompting method (middle part in Figure 2) upon LLMs and multimodal generative models: (1) Text-

to-Image Bridge (T2I-B): we generate the visual imaginative prompt that translates the complex textual instructions (vanilla plan in Figure 3) into explicit scene descriptions (prompt in Figure 3) for text-to-image models. (2) Image-to-Text Bridge (I2T-B): we verbalize the image plan with the image captioning model for generating prompts (red highlighted template in Figure 5) that elicit the revision ability of LLMs with awareness of context. Figure 2 depicts how TIP implements multimodal procedural planning by connecting LLMs and multimodal generative models (Image Caption Model, Text-to-Image Model) with our T2I-B and I2T-B, grounding the image plansin textual context and the text plan in visual context respectively (right part in Figure 2).

### 3.3 Vanilla Text Plan Generation

We first elicit procedural knowledge of LLM to generate vanilla text plan using Zero-shot Chain-of-Thought (Kojima et al., 2022) that does not require heavy human-engineered few-shot examples. Specifically, we leverage InstructGPT (Ouyang et al., 2022) to generate a goal-conditioned step-by-step procedure with the template "*[TEMPLATE]* Task: *[GOAL]*?". *[TEMPLATE]* represents the hand-crafted template to extract the procedural knowledge from LLM. We extend the template "Let's think step by step" (proposed in (Kojima et al., 2022)) as "What's the step-by-step procedure of" for procedural planning. Then we replace the input slot *[GOAL]* with the given task name $T$ (the high-level goal description) as the prompt $P$ to be fed into the LLM. The LLM then outputs goal-conditioned subsequent steps $\mathcal{W} = \{t_1, t_2, ..., t_n\}$ using greedy decoding as our initial textual plan, which is conditioned only on the task name $T$ in zero-shot generation manner.

### 3.4 Textual-Grounded Image Plan Generation with Text-to-Image Bridge

Our Text-to-Image Bridge (T2I-B) in Figure 3 leverages LLM to bridge the gap between the language understanding capabilities of LLM and the ability of language-conditioned image generation in the text-to-image model.

**Imagination Prompt Generation** We encourage LLM to revise the prompt that already processes the physical or temporal meaning residing in the original textual plan. To access this, for each step, we use the prompt $P_{t2i}$ "*[STEP] [T2I-B]*" that concatenates the original generated textual plan at step
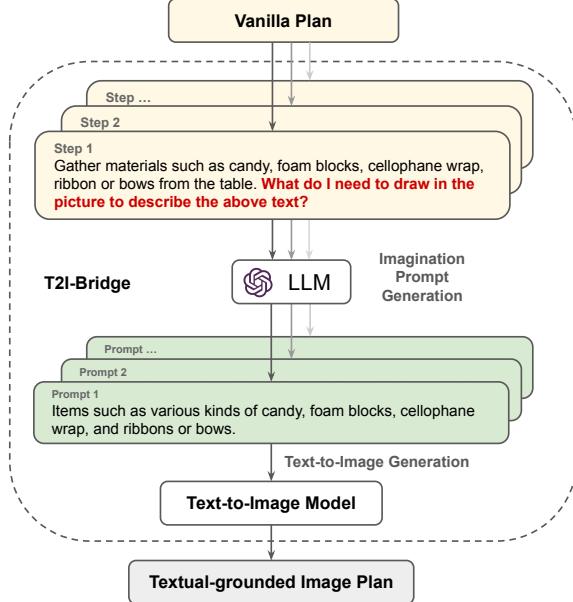


Figure 3: T2I-B elicits visual imagination in LLM to generate explicit scene description (imagination prompt) for text-to-image model conditioned on the vanilla plan.



(a) w. T2I-B   (b) w.o. T2I-B   (c) w. T2I-B   (d) w.o. T2I-B

Figure 4: Text-to-image generation showcases with (a) (b) on "before hitting the waves, read up on ocean safety tips and know the rules of the beach" and (c) (d) on "put down the wine glass" with or without T2I-B.

$i$ and the Text-to-Image Bridge template. *[STEP]* represents one of the subsequent steps generated from LLMs. For *[T2I-B]*, we use the trigger sentence similar to "What do I need to draw in the picture to describe the above text?". With this Text-to-Image Bridge guided prompt $P_{t2i}$, the text-to-image model then generates the textual grounded image at each timestep to compose the final sequence of visual plan $\mathcal{V} = \{v_1, v_2, ..., v_n\}$.

**Text-to-Image Generation** We exploit the Stable Diffusion (Rombach et al., 2022) model to generate RGB images at $512 \times 512$ resolution. Figure 4 provides examples of text-to-image generation with and without our T2I-B. Benefiting from the existing knowledge in LLMs, the text-to-image models are able to generate semantically relevant and high-fidelity images based on the already processed prompt.
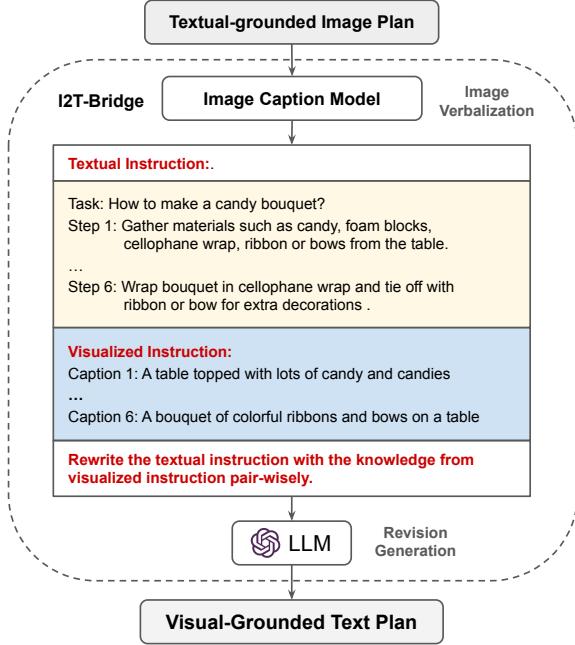
Figure 5: I2T-B injects verbalization of the image plans to foster revision generation of visual-grounded text plans with awareness of multimodal context.

## 3.5 Visual-Grounded Text Plan Generation with Image-to-Text Bridge

To enhance the completeness, alignment, and knowledge exchange between the generated text and image plans, we propose revising the initial text plan based on the textual-grounded image plan.

**Image Verbalization** To complete this, we first need to transfer the visual plan into a natural language format and then inject it into LLM. We implement this by generating captions for each visual plan. Given the image $v$, the captioning model BLIP (Li et al., 2022) generates captions, which transfer the visual knowledge into textual descriptions. For each generated visual plan $v_i$ at each timestep $i$, we generate such pairwise caption with $caption = G(v, Desc)$, where $Desc$ denotes the task description for unified vision and language models, "what does the image describe" in our case. With these image captions, we can further transfer the visual-grounded information into LLM and revise our textual plan.

**Revision Generation** To ground the textual plan in visual context, we use the verbalized description of the visual plan to concatenate with our Image-to-Text Bridge template similar to "Let's revise the procedure using the captions". Concretely, we concatenate the initial textual plan, the captions of the visual plan, and the Image-to-Text Bridge template as the prompt $P_{i2t}$ "Step-by-step Proce-

dure: *[INITIAL]* Captions: *[CAPTION] [I2T-B]*". In this way, we elicit the zero-shot multimodal reasoning ability of LLMs to ground the textual plan in verbalized visual context, as depicted in Figure 5 To this end, our generated multimodal plan is bi-directional grounded by connecting the abilities of LLMs and multimodal generative models.

## 4 Experiments

### 4.1 Datasets

Our datasets are collected and repurposed from WIKIHOW[2] and RECIPEQA(Yagcioglu et al., 2018) due to their temporal relatedness among texts and images. We collect WIKIPLAN by crawling the household "how to" articles from WIKIHOW and then repurpose them into a multimodal procedural planning dataset by formulating the article title as the task name and content as the textual steps, with the pictures as the visual steps. RECIPEQA is a dataset designed for multimodal comprehension of cooking recipes. We collect RECIPEPLAN from this dataset for multimodal procedural planning by sequencing all the given text-image pairs as the text and image plan correspondingly, with the main title as the task name. We conduct zero-shot experiments on $1,000$ distinct, randomly sampled tasks from each dataset. Please refer to Appendix C for more details on the datasets.

### 4.2 Evaluation Metrics

We conduct head-to-head comparisons using Amazon Mechanical Turk (AMT) platform (details can be found in Appendix D.1) on four aspects: (1) `Textual Informativenss`: the text plans contain the necessary information to complete the task, (2) `Visual Informativeness`: the image plans contain the necessary information to complete the task, (3) `Temporal Coherence`: the multimodal plans meet the temporal commonsense requirements, such as the order in which the steps occur, (4) `Planning Accuracy`: whether referring to the multimodal plans can successfully assist task completion. In addition, we measure semantic relevance between predicted text plans and reference text plans using Word Mover's Distance (WMD) (Kusner et al., 2015), Sentence-BERT (S-BERT) (Reimers and Gurevych, 2019), ROUGE-L (Lin, 2004), and METEOR (Banerjee and Lavie, 2005). We measure FID score (Heusel et al., 2017) and CLIP score (Hessel et al., 2021; Radford et al.,

---

[2]https://www.wikihow.com

| Dataset | Ours vs. Model | Textual-Informativeness | | | Visual-Informativeness | | | Temporal Coherence | | | Plan Accuracy | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Win(↑) | Tie | Lose(↓) | Win(↑) | Tie | Lose(↓) | Win(↑) | Tie | Lose(↓) | Win(↑) | Tie | Lose(↓) |
| **WIKIPLAN** | Image Ref + OFA-Caption | **63.34** | 18.38 | 18.27 | **60.63** | 20.45 | 18.92 | **61.95** | 21.03 | 17.02 | **61.99** | 19.40 | 18.61 |
| | Image Ref + BLIP-Caption | **62.70** | 18.70 | 18.60 | **61.26** | 21.18 | 17.56 | **62.22** | 20.78 | 17.00 | **62.29** | 18.28 | 19.43 |
| | Text Ref + DALLE | **62.61** | 20.34 | 17.06 | **59.88** | 22.38 | 17.74 | **60.53** | 22.08 | 17.40 | **61.19** | 22.07 | 16.74 |
| | Text Ref + Stable-Diffusion | **62.58** | 19.82 | 17.60 | **60.25** | 21.16 | 18.58 | **60.68** | 22.38 | 16.94 | **61.73** | 20.56 | 17.72 |
| | Text-Davinci-002 + Stable-Diffusion | **60.68** | 21.56 | 17.76 | **59.90** | 20.41 | 19.70 | **60.22** | 22.99 | 16.79 | **60.41** | 21.53 | 18.06 |
| | Text-Davinci-003 + Stable-Diffusion | **62.32** | 19.82 | 17.86 | **60.29** | 20.85 | 18.85 | **61.10** | 22.17 | 16.73 | **61.48** | 20.29 | 18.23 |
| **RECIPEPLAN** | Image Ref + OFA-Caption | **64.51** | 18.29 | 17.20 | **62.39** | 20.18 | 17.43 | **62.74** | 20.40 | 16.86 | **63.66** | 19.19 | 17.15 |
| | Image Ref + BLIP-Caption | **64.81** | 18.58 | 16.61 | **62.29** | 19.60 | 18.11 | **62.70** | 20.72 | 16.58 | **62.90** | 19.08 | 18.02 |
| | Text Ref + DALLE | **61.16** | 20.15 | 18.69 | **59.60** | 20.60 | 19.80 | **60.04** | 20.48 | 19.48 | **62.11** | 19.21 | 18.68 |
| | Text Ref + Stable-Diffusion | **61.31** | 19.81 | 18.87 | **60.49** | 20.37 | 19.14 | **60.37** | 20.33 | 19.31 | **62.38** | 18.81 | 18.81 |
| | Text-Davinci-002 + Stable-Diffusion | **62.50** | 19.33 | 18.17 | **60.59** | 18.12 | 21.29 | **61.24** | 21.13 | 17.63 | **62.30** | 17.38 | 20.31 |
| | Text-Davinci-003 + Stable-Diffusion | **62.65** | 19.26 | 18.09 | **61.10** | 20.00 | 18.90 | **61.46** | 20.60 | 17.94 | **62.85** | 18.75 | 18.40 |

Table 1: Percentages of multimodal procedural planning results of TIP that are better than, tied with, or worse than baselines, on randomly sampled 200 distinct tasks from each dataset.

| Dataset | Model | Text Plan | | | | Image Plan | | Multimodality Plan | | | Step Length |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | WMD | S-BERT | ROUGE-L | METEOR | FID↓ | CLIP↑ | Cap-S | Text-S | ALL-S | Avg. |
| **WIKIPLAN** | Image Ref + BLIP-Caption | 0.78 | 0.35 | 0.06 | 0.04 | - | 0.71 | <u>0.36</u> | 0.41 | <u>0.39</u> | 8.26 |
| | Image Ref + OFA-Caption | 0.86 | 0.27 | 0.07 | 0.06 | - | 0.71 | 0.27 | <u>0.48</u> | 0.37 | 8.26 |
| | Text Ref + DALLE | 0.68 | <u>0.76</u> | <u>0.28</u> | <u>0.12</u> | 47.39 | 0.74 | 0.33 | 0.26 | 0.29 | 8.26 |
| | Text Ref + Stable-Diffusion | 0.68 | 0.76 | 0.28 | 0.12 | 56.64 | 0.73 | 0.34 | 0.26 | 0.30 | 8.26 |
| | Text-Davinci-002 + Stable-Diffusion | <u>0.87</u> | 0.65 | 0.10 | 0.06 | 61.17 | 0.50 | 0.33 | 0.25 | 0.28 | 4.70 |
| | Text-Davinci-003 + Stable-Diffusion | 0.86 | 0.67 | 0.11 | 0.08 | 57.87 | 0.70 | 0.33 | 0.27 | 0.30 | 6.68 |
| | **TIP (Ours)** | **0.90** | **0.67** | **0.12** | **0.09** | **48.82** | <u>**0.78**</u> | **0.34** | **0.28** | **0.31** | 6.75 |
| **RECIPEPLAN** | Image Ref + BLIP-Caption | 0.77 | 0.37 | 0.08 | 0.05 | - | 0.64 | 0.42 | <u>0.56</u> | <u>0.49</u> | 6.93 |
| | Image Ref + OFA-Caption | 0.82 | 0.40 | 0.09 | 0.10 | - | 0.64 | 0.43 | 0.48 | 0.46 | 6.93 |
| | Text Ref + DALLE | 0.21 | 0.59 | 0.10 | 0.09 | <u>53.55</u> | 0.63 | 0.46 | 0.40 | 0.43 | 6.93 |
| | Text Ref + Stable-Diffusion | 0.21 | 0.59 | 0.10 | 0.09 | 54.58 | 0.61 | <u>0.48</u> | 0.40 | 0.44 | 6.93 |
| | Text-Davinci-002 + Stable-Diffusion | 0.84 | 0.63 | 0.11 | 0.10 | 60.11 | 0.49 | 0.44 | 0.33 | 0.38 | 5.17 |
| | Text-Davinci-003 + Stable-Diffusion | 0.85 | 0.68 | 0.12 | 0.13 | 60.07 | 0.73 | 0.42 | 0.35 | 0.38 | 6.82 |
| | **TIP (Ours)** | <u>**0.86**</u> | <u>**0.68**</u> | <u>**0.13**</u> | <u>**0.14**</u> | **58.68** | <u>**0.79**</u> | **0.43** | <u>**0.36**</u> | <u>**0.40**</u> | 6.94 |

Table 2: Automatic evaluations on $2,000$ distinct tasks from WIKIPLAN and RECIPEPLAN. Image Ref and Text Ref baselines use image and text title references from the dataset. Our TIP uses Text-Davinci-003 and Stable-Diffusion as the LLM and T2I model. We <u>underline</u> and **bold** highest score of models with and without reference baselines.

2021) of image plans. We compute S-BERT between the captions of the predicted image plan and the reference text plan as the Caption-Sentence-BERT score (Cap-S), and between the predicted text plan and the reference text plan as the Text-Sentence-BERT score (Text-S). We then average these two scores to obtain the All-Sentence-BERT score (ALL-S) for multimodal plans. Evaluations are conducted at a procedure level.

### 4.3 Baselines

The key ingredient of our proposed method TIP is that LLMs and multimodal generative models will collaboratively generate multimodal procedural plans benefiting from our designed dual bridges: Text-to-Image Bridge and Image-to-Text Bridge. We compare TIP with the following baselines: (1) ImageRef + OFA/BLIP-Caption: use image plans references directly from the dataset, and generate text plans using image caption models over references (2) TextRef + DALLE/Stable-Diffusion:

use title references from dataset as text plans, and use text-to-image models to generate text plans (3) Text-Davinci-002/003 + Stable-Diffusion: separately generate text and image plans using LLMs and text-to-image models (4) Text-Davinci-003 (Step-based) + Stable-Diffusion: instead of generating the plan at the procedure level, it generates each step iteratively by prompting generated plans from the history sequences to LLMs.

### 4.4 Quantitative Analysis

**Human Evaluation Results** We conduct Win-Tie-Lose Comparison between TIP and the baselines over WIKIPLAN and RECIPEPLAN. Averaged results from 200 tasks rated by 3 crowdsourcing per example are reported in Table 1. Across four aspects, TIP receives consistently higher preferences, outperforming the baselines over the winning ratio by over $60\%$. In terms of textual informativeness, the unimodal baselines (Image Ref + OFA-Caption and Image Ref + BLIP-Caption) is slightly

| Text-to-Image Bridge Template | Alignment | | Image-to-Text Bridge Template | Alignment | |
|---|---|---|---|---|---|
| | WIKIPLAN | RECIPEPLAN | | WIKIPLAN | RECIPEPLAN |
| What do I need to draw in the picture to describe the above text? | **0.9625** | **0.9595** | Rewrite the textual instruction with the knowledge from visualized instruction pair-wisely. | 0.7644 | 0.6945 |
| What do you see in the figure? | 0.9366 | 0.9397 | Based on the visual caption, can you revise the step-by-step procedure according to the paired captions? | **0.8011** | 0.6205 |
| Describe what the picture corresponding to the text should have. | 0.9070 | 0.9181 | Revise each step according to the visual imagination. | 0.6921 | 0.7329 |
| Let's think about what we need to visualize to present the above idea. | 0.8986 | 0.8941 | Let's revise the procedure using the captions. | 0.6155 | **0.7691** |
| Describe something irrelevant to the above text. | 0.5598 | 0.5325 | What's the procedure that disobey the captions? | 0.5079 | 0.5902 |
| What do you usually draw? | 0.5350 | 0.4826 | Provide an interesting procedure to be irrelevant with the captions. | 0.1519 | 0.163 |

Table 3: Robustness check of various templates used in both Text-to-Image Bridge and Image-to-Text Bridge over WIKIPLAN and RECIPEPLAN dataset. The underlined templates are misleading examples. Our Text-Image Prompting model chooses the template with averaged best multimodal alignment, highlighted in purple.



(a) Multimodal procedural plan generated by baseline Text-Davinci-003 + Stable-Diffusion.



(b) Multimodal procedural plan generated by our Text-Image Prompting (TIP).

Figure 6: Improved grounding in textual and visual context are highlighted in pink and green respectively. Red texts indicate reasoning of physical action in image plan generation.

worse than the unimodal text reference based baseline (Text Ref + Stable-Diffusion and Text Ref + DALLE) and multimodal baselines (Text-Davinci-003 + Stable-Diffusion and Text-Davinci-002 + Stable-Diffusion). This is mainly due to the other baselines either direclty leverage the textual information from the reference or the rich text-based knowledge in LLMs. In terms of visual informativeness, the multimodal baselines (Text-Davinci-003 + Stable-Diffusion and Text-Davinci-002 + Stable-Diffusion) can not achieve on par results

with textual reference-based baseline. We hypothesize this is due to the lack of visual knowledge injected into LLMs. The large performance gain of TIP over multimodal baselines (Text-Davinci-003 + Stable-Diffusion and Text-Davinci-002 + Stable-Diffusion) that simply combine the knowledge from LLMs and multimodal generative models imply the importance of grounding our multimodal plans in a multimodal context.

**Automatic Evaluation Results** In Table 2, TIP achieves consistent improvement over baselines

(without Ref), and even surpasses the baselines using reference from the dataset on RECIPEPLAN. This further confirms our superiority in generating multimodal plans with semantic correctness and alignment. Notice that Text Ref baselines directly use the title from the dataset, which is a summarized version of the main content (golden reference used in automatic evaluations).

**Template Robustness** In Table 3, we compare various similar templates for T2I-B and I2T-B against misleading templates. The Alignment is measured with CLIP (Radford et al., 2021) to capture the similarity between given text/image and conditionally generated image/text. The poor alignment of misleading templates and similar alignment of various bridge templates prove the robustness of the template choice in the experiments.

## 4.5 Qualitative Analysis

**Multimodal Grounding** In Figure 6, we compare the performance of TIP to baselines in multimodal procedural planning. TIP generate image plans that are grounded in the textual context. With the help of LLMs reasoning in the temporal dimension, we transfer this ability to image generation, conditioning on the revised prompts of LLMs. This allows digestion of the temporal and complex reasoning present in the text plan and directly indicates what needs to be depicted in the image. The highlighted steps of image plans correctly visualize the scene described in the textual context. For example, at Step 2, instead of only showing the vegetables, ours show an image of a person washing the produce thoroughly. TIP also generate text plans that are better grounded in the image plan. The text plan correctly refers to the objects in visual input, such as "liquid form" and "blended mixture", and also complements the visual context, such as "extract juice from it". Our results indicate the potential for uncovering multimodal reasoning capabilities in LLMs, even though they are primarily used for language reasoning. We provide more comparisons on multimodal procedural planning in Appendix E.1.

## 4.6 Ablations

**Bridge Effect** We report the performance drop of TIP without Image-to-Text Bridge on average textual metrics in Table 4, indicating that the text plan without condition on visual information is vulnerable in text-only planning quality. Then we ablate the text-to-image models in Table 5. With obvious improvement over both FID score and Alignment,

| Plan w.o. I2T-B | WIKIPLAN | RECIPEPLAN |
|---|---|---|
| | Avg. Textual | Avg. Textual |
| Imagination Prompt | 0.341 (-18.4%) | 0.363 (-14.1%) |
| Image Verbalization | 0.261 (-37.5%) | 0.273 (-35.4%) |

Table 4: Ablation of Image-to-Text Bridge on TIP over text plan generation.

| Model | WIKIPLAN | | RECIPEPLAN | |
|---|---|---|---|---|
| | FID ↓ | Align ↑ | FID ↓ | Align ↑ |
| DALLE | 119.03 | 0.77 | 83.27 | 0.64 |
| + Text-to-Image Bridge | **117.02** | **0.79** | **67.64** | **0.78** |
| Stable Diffusion | 129.13 | 0.74 | 88.17 | 0.62 |
| + Text-to-Image Bridge | **119.74** | **0.78** | **84.37** | **0.78** |

Table 5: Ablation of Text-to-Image Bridge with Text Ref + Stable-Diffusion on single-step image plan generation.
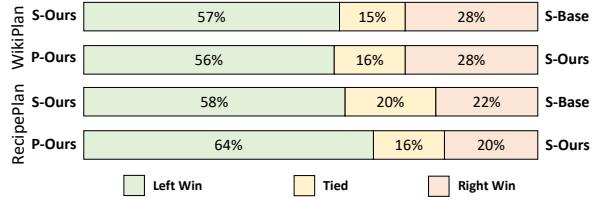


Figure 7: Step-based (S) vs. Procedure-based (P) Win-Tie-Lose over `Plan Accuracy`.

we show that the T2I-B is essential to generate textual-grounded image plans.

**Step-based or Procedure-based** We explore our procedure-based method (P-Ours) against the step-based TIP (S-Ours) and step-based Text-Davinci-003 + Stable-Diffusion (S-Base). We report the head-to-head comparison results on `Plan Accuracy` in Figure 7. The procedure-based method achieves $60\%$ win rate over the step-based TIP. We observe this is partially due to the instinct of LLMs to repeat input texts and is less clear to understand the full intent of generation expectation. Thus the procedure-based method usually achieves better planning quality at the very beginning. We also show that our devised T2I-B and I2T-B achieve consistent improvement on step-based mehtod with averaged winning rate of $58\%$ (S-Ours vs. S-Base).

## 5 Conclusion and Future Work

We introduce the Multimodal Procedural Planning task that aims to generate goal-conditioned text and image subsequences and benchmark models' performance with our curated testbed WIKIPLAN

and RECIPEPLAN. We propose Text-Image Prompt (TIP), a dual-modality prompting framework, that connects LLMs with multimodal generative models to enable plausible multimodal procedural plan generation. Our evaluation benchmark is limited in that no perfect metrics exist to quantify the quality of text-image plans. We hope our work shed light on research into uncovering this limitless capability of multimodal procedural planning driven by uniform automatic metrics.

## Limitations

Relying on the LLMs to reason over complex text for text-to-image models though improving the quality, still remains a large gap with human performance. This is mainly restricted by the pre-training gap between LLMs and text-to-image models. To solve this, further work should explore the finetuning stage that how to inject this language reasoning into the multimodal generation models.

In addition to the model-side limitations, the dataset is limited in that not all the possible multimodal plans are provided and their quality is hard to validate. Due to the lack of perfect metrics in evaluating the text-image sequences, the research in multimodal procedural planning maybe difficult to scale up. Future work should explore this promising direction and furthermore lead LLMs and T2I models better multimodal procedural planners.

## Ethics Statement

We acknowledge that our research utilizes resourceful knowledge in large-scale pre-trained models, which have the potential to bias to a certain cultural background. For example, the task from RECIPEPLAN and WIKIPLAN that involve food preparation may have different procedures depending on different individuals' eating habits. We encourage future studies that uncover the multimodal procedural planning ability with consideration of personalized decision makings.

The data annotation part of the project is classified as exempt by Human Subject Committee via IRB protocols. The hourly wage paid to participants is estimated at $12, which is higher than the federal minimum wage. We manually ensure no personal information is collected and no offensive content is presented during human evaluations.

## References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*.

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3674–3683.

Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *IEEvaluation@ACL*.

Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. 2022. Do as i can, not as i say: Grounding language in robotic affordances. In *6th Annual Conference on Robot Learning*.

Tim Brooks, Aleksander Holynski, and Alexei A Efros. 2022. Instructpix2pix: Learning to follow image editing instructions. *arXiv preprint arXiv:2211.09800*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:1877–1901.

Rémi Calizzano, Malte Ostendorff, and Georg Rehm. 2021. Ordering sentences and paragraphs with pre-trained encoder-decoder transformers and pointer ensembles. In *Proceedings of the 21st ACM Symposium on Document Engineering*, pages 1–9.

Tuhin Chakrabarty, Arkady Saakyan, Olivia Winn, Artemis Panagopoulou, Yue Yang, Marianna Apidianaki, and Smaranda Muresan. 2022. I spy a metaphor: Large language models and diffusion models co-create visual metaphors. *OpenReview Preprint*. Preprint under review.

Chien-Yi Chang, De-An Huang, Danfei Xu, Ehsan Adeli, Li Fei-Fei, and Juan Carlos Niebles. 2020. Procedure planning in instructional videos. In *European Conference on Computer Vision*, pages 334–350. Springer.

Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 2016. Neural sentence ordering. *arXiv preprint arXiv:1607.06952*.

Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. 2021. Unifying vision-and-language tasks via text generation. In *International Conference on Machine Learning*, pages 1931–1942. PMLR.

Baiyun Cui, Yingming Li, Ming Chen, and Zhongfei Zhang. 2018. Deep attentive sentence ordering network. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4340–4349.

Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. 2021. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835.

Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. 2022. Make-a-scene: Scene-based text-to-image generation with human priors. *arXiv preprint arXiv:2203.13131*.

Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *NIPS*.

Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. 2022. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10696–10706.

Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. CLIPScore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6629–6640, Red Hook, NY, USA. Curran Associates Inc.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851.

Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *arXiv preprint arXiv:2201.07207*.

Vladimir Ilievski, Claudiu Musat, Andreea Hossmann, and Michael Baeriswyl. 2018. Goal-oriented chatbot dialog management bootstrapping with transfer learning. In *IJCAI*.

Peter Jansen. 2020. Visually-grounded planning without vision: Language models infer detailed plans from high-level instructions. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4412–4417.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners.

Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. 2015. From word embeddings to document distances. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 957–966, Lille, France. PMLR.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *ICML*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Vivian Liu and Lydia B Chilton. 2022. Design guidelines for prompt engineering text-to-image generative models. In *CHI Conference on Human Factors in Computing Systems*, pages 1–23.

Xiao Liu, Da Yin, Yansong Feng, and Dongyan Zhao. 2022. Things not written in text: Exploring spatial commonsense from visual signals. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2365–2376.

Yujie Lu, Weixi Feng, Wanrong Zhu, Wenda Xu, Xin Eric Wang, Miguel Eckstein, and William Yang Wang. 2022a. Neuro-symbolic causal language planning with commonsense prompting. *arXiv preprint arXiv:2206.02928*.

Yujie Lu, Wanrong Zhu, Xin Wang, Miguel Eckstein, and William Yang Wang. 2022b. Imagination-augmented natural language understanding. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pages 4392–4402, Seattle, United States. Association for Computational Linguistics.

Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*.

Byungkook Oh, Seungmin Seo, Cheolheon Shin, Eunju Jo, and Kyong-Ho Lee. 2019. Topic-guided coherence modeling for sentence ordering by preserving global and local information. In *Proceedings of*

the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2273–2283.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. 2022. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155.

Liang Qiu, Yuan Liang, Yizhou Zhao, Pan Lu, Baolin Peng, Zhou Yu, Ying Nian Wu, and Song-Chun Zhu. 2021. Socaog: Incremental graph parsing for social relation inference in dialogues. In *The 59th Annual Meeting of the Association for Computational Linguistics (ACL)*.

Liang Qiu, Yizhou Zhao, Yuan Liang, Pan Lu, Weiyan Shi, Zhou Yu, and Song-Chun Zhu. 2022. Towards socially intelligent agents with mental state transition and human utility. In *The 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*.

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*.

Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *ArXiv*, abs/1908.10084.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695.

Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*.

Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR.

Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. 2022. Llm-planner: Few-shot grounded planning for embodied agents with large language models. *arXiv preprint arXiv:2212.04088*.

Yang Song and Stefano Ermon. 2019. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32.

Yixuan Su, Tian Lan, Yahui Liu, Fangyu Liu, Dani Yogatama, Yan Wang, Lingpeng Kong, and Nigel Collier. 2022. Language models can see: Plugging visual controls in text generation. *arXiv preprint arXiv:2205.02655*.

Ming Tao, Hao Tang, Songsong Wu, Nicu Sebe, Xiao-Yuan Jing, Fei Wu, and Bingkun Bao. 2020. Df-gan: Deep fusion generative adversarial networks for text-to-image synthesis. *arXiv preprint arXiv:2008.05865*.

Stefanie Tellex, Thomas Kollar, Steven Dickerson, Matthew Walter, Ashis Banerjee, Seth Teller, and Nicholas Roy. 2011. Understanding natural language commands for robotic navigation and mobile manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1507–1514.

Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. 2021. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212.

Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems*, 30.

Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022a. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. *CoRR*, abs/2202.03052.

Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022b. Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. *arXiv preprint arXiv:2202.03052*.

Qingyun Wang, Manling Li, Hou Pong Chan, Lifu Huang, Julia Hockenmaier, Girish Chowdhary, and Heng Ji. 2022c. Multimedia generative script learning for task planning. *arXiv preprint arXiv:2208.12306*.

Zhenhailong Wang, Manling Li, Ruochen Xu, Luowei Zhou, Jie Lei, Xudong Lin, Shuohang Wang, Ziyi Yang, Chenguang Zhu, Derek Hoiem, et al. 2022d. Language models with image descriptors are strong few-shot video-language learners. *arXiv preprint arXiv:2205.10747*.

Zijie J Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. 2022e. Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models. *arXiv preprint arXiv:2210.14896*.

Te-Lin Wu, Alex Spangher, Pegah Alipoormolabashi, Marjorie Freedman, Ralph Weischedel, and Nanyun Peng. 2022. Understanding multimodal procedural knowledge by sequencing multimodal instructional manuals. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4525–4542.

Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. 2018. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1316–1324.

Semih Yagcioglu, Aykut Erdem, Erkut Erdem, and Nazli Ikizler-Cinbis. 2018. Recipeqa: A challenge dataset for multimodal comprehension of cooking recipes. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1358–1368.

Shiquan Yang, Rui Zhang, Sarah Erfani, and Jey Han Lau. 2022a. An interpretable neuro-symbolic reasoning framework for task-oriented dialogue generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4918–4935.

Yue Yang, Wenlin Yao, Hongming Zhang, Xiaoyang Wang, Dong Yu, and Jianshu Chen. 2022b. Z-lavi: Zero-shot language solver fueled by visual imagination. *arXiv preprint arXiv:2210.12261*.

Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. 2022c. An empirical study of gpt-3 for few-shot knowledge-based vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 3081–3089.

Hui Ye, Xiulong Yang, Martin Takac, Rajshekhar Sunderraman, and Shihao Ji. 2021. Improving text-to-image synthesis using contrastive learning. *arXiv preprint arXiv:2107.02423*.

Andy Zeng, Adrian Wong, Stefan Welker, Krzysztof Choromanski, Federico Tombari, Aveek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, et al. 2022. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*.

Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. 2021a. Cross-modal contrastive learning for text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 833–842.

Li Zhang, Qing Lyu, and Chris Callison-Burch. 2020. Reasoning about goals, steps, and temporal ordering with wikihow. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4630–4639.

Lisai Zhang, Qingcai Chen, Joanna Siebert, and Buzhou Tang. 2021b. Semi-supervised visual feature integration for language models through sentence visualization. In *Proceedings of the 2021 International Conference on Multimodal Interaction*, pages 682–686.

Yizhou Zhao, Liang Qiu, Wensi Ai, Pan Lu, and Song-Chun Zhu. 2022. Triangular character animation sampling with motion, emotion, and relation. *arXiv preprint arXiv:2203.04930*.

Minfeng Zhu, Pingbo Pan, Wei Chen, and Yi Yang. 2019. Dm-gan: Dynamic memory generative adversarial networks for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5802–5810.

Wanrong Zhu, An Yan, Yujie Lu, Wenda Xu, Xin Eric Wang, Miguel Eckstein, and William Yang Wang. 2022a. Visualize before you write: Imagination-guided open-ended text generation. *arXiv preprint arXiv:2210.03765*.

Xizhou Zhu, Jinguo Zhu, Hao Li, Xiaoshi Wu, Hongsheng Li, Xiaohua Wang, and Jifeng Dai. 2022b. Uni-perceiver: Pre-training unified architecture for generic perception for zero-shot and few-shot tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16804–16815.

## A Background

A line of work in unimodal procedural planning studies sorting a series of unordered texts or events (Chen et al., 2016; Cui et al., 2018; Oh et al., 2019; Calizzano et al., 2021; Wu et al., 2022). Other work explores generating subsequent steps given a target goal, e.g., Lu et al. (2022a) aim at generating a sequence of plans to complete the high-level task.

Text-to-image generation is a task that synthesizes images from text prompts. DALL-E 2 (Ramesh et al., 2022) and Stable Diffusion (Rombach et al., 2022) are state-of-the-art text-to-image models developed on top of diffusion models conditioned on input texts. Some early work in text-to-image models trains generative adversarial networks (GANs) (Goodfellow et al., 2014) on image captioning datasets (Xu et al., 2018; Zhu et al., 2019; Tao et al., 2020; Zhang et al., 2021a; Ye et al., 2021) to generate images conditioned on textual descriptions. Other work follows the VQ-VAE (Van Den Oord et al., 2017) framework and trains autoregressive transformers that take both the text and image as sequences of tokens (Ramesh et al., 2021; Ding et al., 2021; Gafni et al., 2022). However, these methods are struggling to generate photorealistic images. Motivated by the remarkable progress of diffusion models in generating images with fidelity (Sohl-Dickstein et al., 2015; Song and Ermon, 2019; Ho et al., 2020), recent work has applied them to text-to-image generation with auxiliary text encoders (Rombach et al., 2022; Nichol et al., 2021; Gu et al., 2022; Ramesh et al., 2022; Saharia et al., 2022). Wang et al. (2022e) propose the first large-scale text-to-image prompt dataset, DiffusionDB, which enables a new research direction of prompt engineering to construct better prompts. Chakrabarty et al. (2022) use GPT-3 to generate a detailed textual description with rich visual metaphors to prompt the DALL-E 2 model.

Recently, there has been a trend of using large language models (LLMs) like GPT-3 (Brown et al., 2020) to transfer visual knowledge in order to improve their capabilities in downstream natural language processing (NLP) and multimodal tasks. For example, images and videos can be translated into captions which further instruct a language model to generate contextual descriptions (Wang et al., 2022d; Zeng et al., 2022) or answer knowledge-based visual questions (Yang et al., 2022c). Instead of being prompted with textual descriptions,

language models can extend to vision-language settings through text generation controlled by visual features (Cho et al., 2021; Tsimpoukelli et al., 2021; Su et al., 2022; Zhu et al., 2022b; Wang et al., 2022b; Alayrac et al., 2022).

## B Method Details

### B.1 Configurations

The experiments using Text-Davinci and DALLE are conducted with OpenAI API on January 2023. We use BLIP w/ ViT-B and CapFilt-L[3] and OFA-base from huggingface demo[4].

### B.2 Details of Module Outputs

We showcase all the details of the outputs of each module for the example task "How to make a candy bouquet" in Figure 8. In addition, we showcase examples on "How to make peppermint patty pudding shot" and "How to make bartending for beginners simple margarita" in Figure 9.

## C Dataset Details

### C.1 RECIPEPLAN

**Data Repurpose** RECIPEQA was proposed in (Yagcioglu et al., 2018) that provide four tasks (Textual Cloze, Visual Cloze, Visual Ordering, Visual Coherence) for multimodal machine comprehension of cooking recipes. This dataset contains question-answer pairs generated from copyright-free recipes. Each of them is under a license, which is provided in each data JSON file. We collect RECIPEPLAN by repurposing the test dataset from RECIPEQA that relates to the Visual Ordering task as the sequence generation task for the multimodal procedural planning evaluation testbed. We use recipe instructions as textual plan reference and their paired images as visual plan reference.

**Dataset Statistics** We visualize two examples of our repurposed RECIPEPLAN for multimodal procedural planning in Figure 10. We also show the word-cloud distribution of task name and textual plan reference in Figure 12.

### C.2 WIKIPLAN

**Raw Data Collection** To facilitate research on learning to generate procedural planning in a multimodal setting, we have constructed the large-scale WIKIPLAN dataset collected from the WIK-
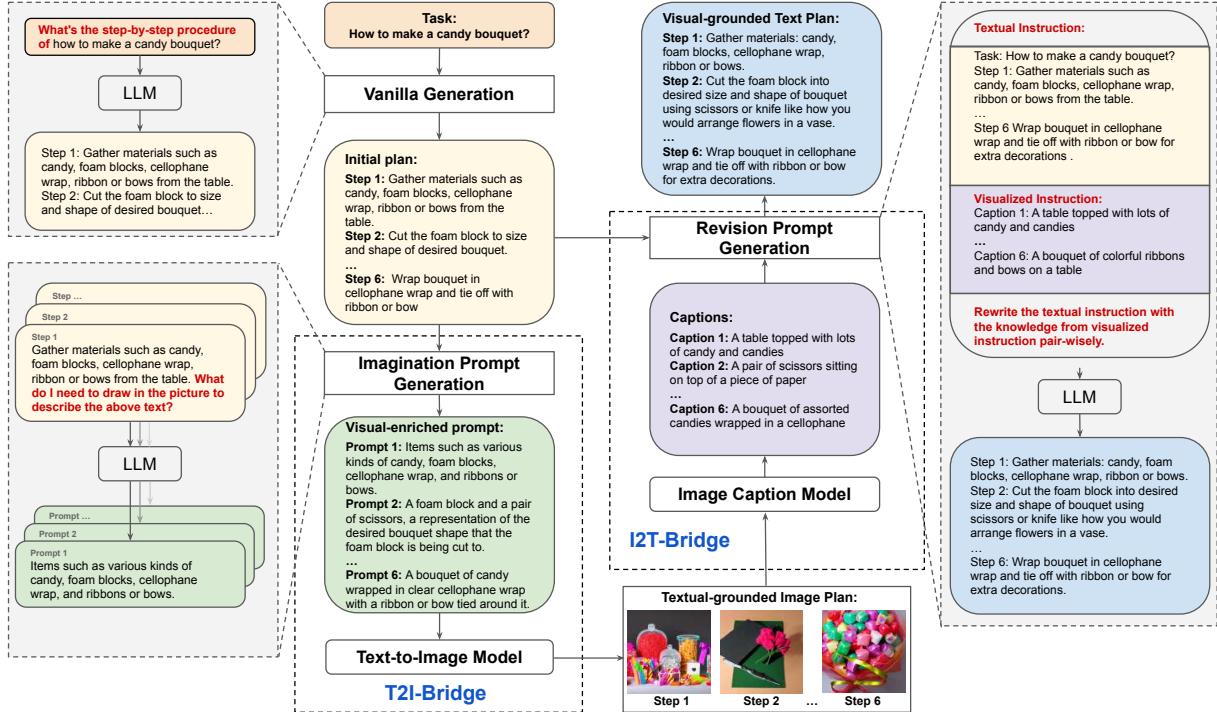
---

[3]https://github.com/salesforce/BLIP
[4]https://huggingface.co/OFA-Sys/ofa-base

**Task:** How to make a candy bouquet?

| Vanilla Text plan | Vanilla Image Plan | Imagined Prompt | Textual-Grounded Image Plan | Verbalization | Visual Grounded Text plan |
|---|---|---|---|---|---|
| **Step 1:** Gather materials: candy, foam blocks, cellophane wrap, ribbon or bows. | | Items such as various kinds of candy, foam blocks, cellophane wrap, and ribbons or bows. | | A table topped with lots of candy and candies | **Step 1:** Gather materials such as candy, foam blocks, cellophane wrap, ribbon or bows from the table. |
| **Step 2:** Cut the foam block to size and shape of desired bouquet. | | A foam block and a pair of scissors, a representation of the desired bouquet shape that the foam block is being cut to. | | A pair of scissors sitting on top of a piece of paper | **Step 2:** Cut the foam block into desired size and shape of bouquet using scissors or knife like how you would arrange flowers in a vase. |
| **Step 3:** Insert wooden skewers into each piece of candy. | | Wooden skewers inserted into each piece of candy, arrows pointing to the skewers to indicate that they are being inserted into the candy. | | A group of toothbrushes sitting on top of a table | **Step 3:** Insert wooden skewers into each piece of candy as if they were lollipops on top of a blue table. |
| **Step 4:** Arrange the pieces of candy onto the foam in a pleasing pattern. | | A bouquet of candy in a pleasing pattern. The candy should be arranged on top of a foam base. The candy should be arranged in an aesthetically pleasing pattern, such as alternating colors, sweet and | | A bouquet of colorful candies in a vase | **Step 4:** Arrange the pieces of candy onto the foam in an interesting pattern like playing with colored rocks on plate and cupcakes and candies for your bouquet decoration ideas. |
| **Step 5:** Secure pieces with hot glue if necessary. | | A bouquet of candy pieces that are being held together with hot glue, a person holding a hot glue gun, to show that the pieces of candy are being secured with hot glue. | | A pile of different colored candies and lollipops | **Step 5:** Secure pieces with hot glue if necessary just like arranging colorful ribbons and bows around the bouquets centerpiece. |
| **Step 6:** Wrap bouquet in cellophane wrap and tie off with ribbon or bow. | | A bouquet of candy wrapped in clear cellophane wrap with a ribbon or bow tied around it. | | A bouquet of assorted candies wrapped in a cellophane | **Step 6:** Wrap bouquet in cellophane wrap and tie off with ribbon or bow for extra decorations. |

(a) Full example outputs.

**What's the step-by-step procedure of** how to make a candy bouquet?

LLM

Step 1: Gather materials such as candy, foam blocks, cellophane wrap, ribbon or bows from the table. Step 2: Cut the foam block to size and shape of desired bouquet…

Step …
Step 2
**Step 1**
Gather materials such as candy, foam blocks, cellophane wrap, ribbon or bows from the table. **What do I need to draw in the picture to describe the above text?**

LLM

Prompt …
Prompt 2
**Prompt 1**
Items such as various kinds of candy, foam blocks, cellophane wrap, and ribbons or bows.

**Task:**
**How to make a candy bouquet?**

**Vanilla Generation**

**Initial plan:**
**Step 1:** Gather materials such as candy, foam blocks, cellophane wrap, ribbon or bows from the table.
**Step 2:** Cut the foam block to size and shape of desired bouquet.
...
**Step 6:** Wrap bouquet in cellophane wrap and tie off with ribbon or bow

**Imagination Prompt Generation**

**Visual-enriched prompt:**
**Prompt 1:** Items such as various kinds of candy, foam blocks, cellophane wrap, and ribbons or bows.
**Prompt 2:** A foam block and a pair of scissors, a representation of the desired bouquet shape that the foam block is being cut to.
...
**Prompt 6:** A bouquet of candy wrapped in clear cellophane wrap with a ribbon or bow tied around it.

**Text-to-Image Model**
**T2I-Bridge**

**Textual-grounded Image Plan:**
Step 1   Step 2   ...   Step 6

**Visual-grounded Text Plan:**
**Step 1:** Gather materials: candy, foam blocks, cellophane wrap, ribbon or bows.
**Step 2:** Cut the foam block into desired size and shape of bouquet using scissors or knife like you would arrange flowers in a vase.
...
**Step 6:** Wrap bouquet in cellophane wrap and tie off with ribbon or bow for extra decorations.

**Revision Prompt Generation**

**Captions:**
**Caption 1:** A table topped with lots of candy and candies
**Caption 2:** A pair of scissors sitting on top of a piece of paper
...
**Caption 6:** A bouquet of assorted candies wrapped in a cellophane

**Image Caption Model**

**I2T-Bridge**

**Textual Instruction:**
Task: How to make a candy bouquet?
Step 1: Gather materials such as candy, foam blocks, cellophane wrap, ribbon or bows from the table.
...
Step 6 Wrap bouquet in cellophane wrap and tie off with ribbon or bow for extra decorations .

**Visualized Instruction:**
Caption 1: A table topped with lots of candy and candies
...
Caption 6: A bouquet of colorful ribbons and bows on a table

**Rewrite the textual instruction with the knowledge from visualized instruction pair-wisely.**

LLM

Step 1: Gather materials: candy, foam blocks, cellophane wrap, ribbon or bows.
Step 2: Cut the foam block into desired size and shape of bouquet using scissors or knife like how you would arrange flowers in a vase.
...
Step 6: Wrap bouquet in cellophane wrap and tie off with ribbon or bow for extra decorations.

(b) Visualization of outputs from each module.

Figure 8: Output Details of each module of TIP for Multimodal Procedural Planning. The T2I-B leverage the complex language comprehension and zero-shot reasoning ability of LLMs to improve text-to-image generation. Reversely, the I2T-B injects visual knowledge via verbalization of the visual plans to generate a visually-grounded and complementary textual plan.

**Task:** How to make peppermint patty pudding shot?

| Vanilla Text plan | Vanilla Image Plan | Imagined Prompt | Textual-Grounded Image Plan | Verbalization | Visual Grounded Text plan |
|---|---|---|---|---|---|
| **Step 1:** Combine 1/2 cup of peppermint schnapps, 1/4 cup of Bailey's Irish cream, and 2 tablespoons of chocolate syrup in a bowl or blender. |  | A bowl or blender filled with 1/2 cup of peppermint schnapps, 1/4 cup of Bailey's Irish cream, and 2 tablespoons of chocolate syrup. |  | Two cups of hot chocolate with spoons on a table | **Step 1:** Place three cups of chocolate pudding on a flat surface and add 1/2 cup of peppermint schnapps, 1/4 cup of Baileys Irish Cream, and 2 tablespoons of chocolate syrup. |
| **Step 2:** Add 4 ounces of cold heavy whipping cream to mixture and blend until smooth. |  | A bowl with the mixture in it, a spoon in the mixture, a measuring cup with 4 ounces of heavy cream, a blender, and a glass with the peppermint patty pudding shot. |  | A table topped with three different types of desserts. | **Step 2:** Using a spoon or blender mix all ingredients until smooth. |
| **Step 3:** Place 4 small glasses on a flat surface and fill with pudding shot mix. |  | Four small glasses on a flat surface, each filled with a pudding shot mix. You should also include any relevant ingredients, such as peppermint extract, liqueur, and cream, to complete the pudding. |  | A table topped with three shot glasses filled with liquid. | **Step 3:** Pour the mixture into four small glasses and top each with whipped cream. |
| **Step 4:** Chill for at least 30 minutes before serving. |  | A picture of a tray of peppermint patty pudding shots in the refrigerator, with a timer counting down the 30 minutes. |  | A tray filled with chocolate cups and candy canes. | **Step 4:** Add a candy cane for garnish on each glass and chill for at least 30 minutes before serving. |

**Task:** How to make bartending for beginners simple margarita?

| Vanilla Text plan | Vanilla Image Plan | Imagined Prompt | Textual-Grounded Image Plan | Verbalization | Visual Grounded Text plan |
|---|---|---|---|---|---|
| **Step 1:** Gather the following ingredients: tequila, triple sec (or Cointreau), lime juice, and simple syrup. |  | A glass of margarita with a lime wedge garnish, and next to it illustrate the ingredients needed to make the margarita - a bottle of tequila, bottle of triple sec or Coint |  | A bottle of orange juice next to limes and lime slices. | **Step 1:** Gather the following ingredients: tequila, triple sec (or Cointreau), lime juice, and simple syrup. |
| **Step 2:** Fill a shaker with ice cubes. |  | A shaker filled with ice cubes. |  | A glass filled with ice and lime slices. | **Step 2:** Fill a shaker with ice cubes. |
| **Step 3:** Pour in 1 ounce of tequila, ½ ounce of triple sec or Cointreau, ¾ ounces of fresh lime juice, and ½ ounce of simple syrup into the shaker filled with ice cubes. |  | A shaker filled with ice cubes, a bottle of tequila, a bottle of triple sec or Cointreau, a bottle of fresh lime juice, and a bottle of simple syrup. |  | A glass of lemonade next to a bottle of lemonade. | **Step 3:** Pour in 1 ounce of tequila, ½ ounce of triple sec or Cointreau, ¾ ounces of fresh lime juice into the shaker filled with ice cubes and shake well until all ingredients are combined and chilled (about 20 seconds). |
| **Step 4:** Shake well until all ingredients are combined and chilled (about 20 seconds). |  | A person shaking a cocktail shaker filled with ice cubes, limes, tequila, and triple sec. You could also add a timer set to 20 seconds to illustrate how long the drink should be shaken. |  | A person pouring a drink into a glass. | **Step 4:** Add half an ounce of simple syrup before straining the mixture into a glass rimmed with salt or sugar-lime mix (optional). |
| **Step 5:** Strain the mixture into a glass rimmed with salt or sugar-lime mix (optional). |  | A glass with a rim lined with either salt or a sugar-lime mix and a pitcher or shaker filled with a mixture of margarita. A spoon or strainer. |  | A glass of water with a lime next to it. | **Step 5:** Garnish with a wedge of lime or orange slice before serving cold margarita in your desired glass! |
| **Step 6:** Garnish with a wedge of lime or orange slice before serving cold margarita to your guests! |  | A glass filled with a cold margarita, a wedge of lime or orange slice on the side of the glass, and possibly some straws. |  | A glass of orange juice with a straw and lime. | **Step 6:** Enjoy your delicious homemade drink! |

Figure 9: Full example of outputs details of TIP for Multimodal Procedural Planning.

IHOW website[5], which is under an Attribution-Noncommercial-Share Alike 3.0 Creative Commons License.. This website provides a wide range of how-to articles related to everyday life topics, which are collaboratively written by its users and reviewed by experts. We crawled each article, collecting the task title, URL, introduction, topics, and steps. Each step includes a brief textual action, a detailed context, and an illustration image. Our raw dataset consists of 30,026 examples across 19 categories and 2,062 topics. We plan to release the

| Task: How to make baileys hot chocolate mason jar. | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Step 1: Gather Your Supplies.** Wide mouth, pint sized mason jars Hot chocolate powder Mini-marshmallows Twine Baileys mini bottles - one for each jar you plan to make. | **Step 2: Wash Your Mason Jars.** Make sure the jars are clean and free of debris. Dry thoroughly before moving on to the next step. | **Step 3: Add Hot Cocoa Mix.** Pour the cocoa mix into the jar. I added 2 bags of powder so it can fill up to a little less than half of the jar. | **Step 4: Marshmallows.** I'm a bit marshmallow fanatic, so I'm friendly with the amount I add. Feel free to add more or less, depending on your preference. Then, close the jar. | **Step 5: Wrap Twine Around the Lid.** Cut a long piece of twine. Firmly wrap the twine several times around the lid, being careful not to let it slip down. Once you have it firmly wrapped around, tie a knot. Leave remaining twine to attach mini bottle. | **Step 6: Tie the Mini Baileys Bottle.** With the remaining twine, tie a simple knot around the mini bottle of Baileys. I would double knot as well, around the mini bottle to be sure it is firm and tightly attached to the mason jar. Cut any excess twine. | **Step 7: This Is What One Bottle Looks Like!** Only a few items needed, to make a fast, unique gift! A person receiving this gift would only need to add hot water or milk and pour in some yummy Baileys to their liking! | **Step 8: Perfect for the Holidays!** This idea is a great gift for any occasion, but it's a great idea for the holidays coming up. You can even add a greeting card or any other embellishment you would like! |

| Task: How to make vegan lavender shortbread cookies. | | | | | | |
|---|---|---|---|---|---|---|
| **Step 1: Prepare the Ingredients.** Ingredients (for 20 cookies): 1/2 cup (115g) margarine at room temperature (you can also use coconut oil); 1 cup (130g) all-purpose white flour; 1/4 cup (30g) powdered sugar; … | **Step 2: Chop the Lavender and Cream the Margarine.** Preheat oven to 350°F / 180°C (gas mark 4) and prepare a cookie tray. Line the tray with parchment paper. In a small bowl sift together the flour and salt. … | **Step 3: Add the Dry Ingredients.** Add the powdered sugar to the large bowl and beat until smooth and creamy. I like to sift my powdered sugar before adding it to the bowl. … | **Step 4: Form a Log.** At this point the mixture should resemble to wet sand. Continue mixing with your fingers until you form a ball of dough. Place the dough between two cling film rectangles and shape it into a log. … | **Step 5: Cut the Cookies.** Remove the log from the freezer and discard the cling film. Flour the working surface. Trim the rounded edges of the log. With superficial slits, divide the log into 5mm-thick cookies. … | **Step 6: Refrigerate and Bake the Cookies.** Place the shortbread cookies on the parchment paper-lined tray making sure you space the cookies at least 5 cm apart. Place the tray in the refrigerator for 15 minutes. … | **Step 7: Enjoy!** Store the cookies in an airtight container at room temperature for up to 1 week. Note: For this recipe I used margarine. You can also use vegan butter or solidified coconut oil, refined or unrefined. … |

Figure 10: Two examples in the RECIPEPLAN dataset.

raw data in the hopes of pre-training models for procedural planning and knowledge reasoning.

**Quality Control** To improve the evaluation of different baselines, we further selected five categories that feature temporal actions and high-fidelity visual contexts: *Food and Entertaining*, *Hobbies and Crafts*, *Home and Garden*, *Pets and Animals*, and *Sports and Fitness*. In order to reflect common tasks in real-life scenarios, articles with fewer than three steps or more than 22 steps were excluded, as well as articles with images of a dimension size of fewer than 400 pixels. Each category was balanced with 200 examples. To further ensure high quality, we conducted a quality control in which well-trained human annotators reviewed the dataset and manually revised the examples if there was any wrong or inappropriate content.

**Dataset Statistics** Finally, our WIKIPLAN dataset consists of 1,000 examples across 5 categories and 370 different topics. Three examples of the dataset are illustrated in Figure 11, which include muddling mint leaves for a cocktail, encouraging a cat to eat, and becoming a better football player. Each

example is composed of a title, introduction, related topics, and a list of detailed steps with visual aids. The word-cloud distributions of the task titles and step text are visualized in Figure 13, which demonstrate the diversity of topics and words.

## D Evaluation Details

### D.1 Crowdsourcing Human Evaluation

We manually ensure no personal information is collected and no offensive content is presented during human evaluations. The hourly wage paid to participants is estimated at $12. And the total amount spent on participant compensation is $1958.

We average the results from 3 annotators for each example. Given the high-level goal (task name) for each assignment, we want the annotators to compare two generated text and image sequences in terms of *Textual-Informativeness*, *Visual-Informativeness*, *Temporal Coherence* and *Plan Accuracy*. Before going to the question, we let the annotators read the instructions below:

**Instruction:** Given the Task (e.g, Task: How to muddle), please compare two sequences of steps

**Task: How to Muddle**

**Introduction:** Muddling is a bartender's technique for releasing flavor from solid cocktail ingredients. The basic idea is as simple as crushing the fruit or herb, but there are devils in the details ready to make your cocktail bitter or unpleasant. Learn how to do it the right way, and you won't have to muddle through on your own intuition.

**Topics:** Spirits and liqueurs, Cocktails

| **Step 1: Choose a gentle muddling tool.** | **Step 2: Put the ingredients in a sturdy glass.** | **Step 3: Press and twist lightly.** | **Step 4: Finish the drink.** |
|---|---|---|---|
| A small, flat wooden tool is ideal, such as the end of a wooden spoon or a French rolling pin (without handles). Plastic or hard rubber tools will also work, but require a delicate touch. Avoid muddlers with teeth, since these tear up the leaves too much. A tough herb such as rosemary needs more breaking down. Follow the instructions for fruit instead. For instance, try using the wide, flat end of a wooden spoon. That will crush whatever you're muddling, without destroying it. | Do not muddle in a delicate glass that may chip or break when struck. If the cocktail also calls for fruit, cucumber, or other non-herb ingredients, muddle them separately for best results. Granulated sugar tears into the ingredients as you muddle. This can be overkill for soft herbs, so add it with the fruit instead, or dissolve it in a few drops of water and add it to the cocktail separately. | Mint, basil, and other soft leaves release bitter flavors when crushed or torn apart. Just press down gently with the blunt tool while you twist your wrist, then release. Do this two or three times. Use your dominant hand to muddle while you hold the glass steady with your other hand. You only need to lightly bruise the leaves to release the oils and aroma. | The leaves are ready when lightly bruised, but still intact. You should be able to smell the herb, since the goal of muddling is to release flavorful, aromatic oils. You can leave the herbs in the final cocktail for presentation or strain them out, as desired. |

**Task: How to Encourage Your Cat to Eat**

**Introduction:** Cats are finicky creatures, and they may go off food when they are sick or getting older, for instance. If your cat suddenly stops eating, you should definitely take it to the vet for a checkup, as it may have an underlying problem. Otherwise, you can work on the cat, environment, and food to encourage your cat to eat, as well as use a few tricks when switching foods.

**Topics:** Cats, Feeding cats

| **Step 1: Feed the cat alone.** | **Step 2: Give smaller meals.** | **Step 3: Offer comfort and attention.** | **Step 4: Wipe the cat's nose.** | **Step 5: Hand feed the cat.** |
|---|---|---|---|---|
| If there are other animals in the house, the cat may feel uncomfortable eating, as it tries to compete for food. If you section off a room where you can feed just that cat, it can increase the chances of the cat eating. | Smaller meals throughout the day may seem counter-intuitive, but it can actually encourage your cat to eat. If your cat doesn't have a large appetite, it may be overwhelmed by a large bowl of food. | Some cats want attention when they're eating. Try stroking your cat and talking soothingly to it while it's eating. You can also try a bit of petting if your cat is near the bowl but not eating. This tactic won't work with every cat, though, so if your cat seems disturbed by the attention, leave it alone. | If the cat has been sick, it may not be smelling very well. Wiping the cat's nose and trying to remove discharge may help it smell better. In turn, it may be more interested in food because it can smell it. | A cat who hasn't been interested in food may be more inclined to eat if you hand feed it. For canned food, you can put a little on your finger and offer it to the cat. For dry food, place a bit in your hand, and hold it out for the cat to eat. |

**Task: How to Be a Better Softball Player**

**Introduction:** Softball is a fun and thrilling sport that requires a lot of hard work and focus. If you feel like you aren't living up to your full potential as a softball player, consider re-evaluating your techniques as you bat, throw, and catch. Additionally, take a moment to think about the demands of your current position, and how you can adjust your practice schedule to hone your skills on the field. With enough focus and dedication, you'll start to see your softball skills improve!

**Topics:** Team sports, Softball

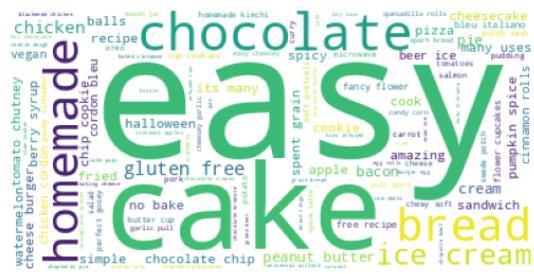| **Step 1: Keep your elbow higher so you can throw farther.** | **Step 2: Thrust your palm forward before letting go with your fingers.** | **Step 3: Practice throwing balls heavier than 6 ounces (0.17 kg) to improve your strength.** | **Step 4: Move your arm in a fluid motion to give the pitch momentum.** | **Step 5: Eye your target before you release the ball.** |
|---|---|---|---|---|
| Pitch the ball with your arm positioned above your shoulder. While the throwing motion will eventually bring your arm to waist or hip level, try to keep your elbow at a higher starting level. When you create the round throwing motion with your arm and elbow, a higher starting point will create more momentum for the ball. Make sure that you're using a fielding glove when you practice, and not a fielding mitt. | Aim to hold the ball in a four-seam grip, where your pointer, middle, and ring finger are positioned behind the edge of the seam you're choosing to hold. When you're holding the ball, don't rest it against your palm. Instead, use palm as a source of momentum as you start to throw, then lift up your fingers quickly to release the ball. One of the most important elements of this type of throw is timing. Practice releasing your fingers at the exact moment your palm hits so that the throw can get more speed. | Build muscle in your arms by practicing with softballs that are heavier than what you plan on using in an actual game. Keep your throwing technique consistent, even if your arm feels a bit more tired. When you play with a normal-sized ball, your arm will feel much more adjusted to the throw. | Continue throwing the ball with your usual technique, using a circular motion as your arm swings over your shoulder to waist or hip level. As you swing downwards, don't slow your movement; instead, keep your arm moving swiftly when you release the ball. If you do this, you're more likely to release a more powerful pitch. If you slow down your arm at the end of the throw, you can lose pitch speed. | Whether you're pitching or fielding, make firm eye contact with where you want to throw the ball. Don't let your gaze waver; even if your aim is not always consistent, you're much more likely to have success if you throw with a target in mind. As you throw, keep your shoulders and hips angled in the direction of this target. Be sure to square your shoulders in a specific direction whenever you throw the ball. Your shoulders help serve as an invisible guide to which direction the ball goes. |

Figure 11: Three examples in our curated WIKIPLAN dataset.

Sequence 1 and Sequence 2, and determine which sequence is better in terms of four aspects:

- Textual-Informativeness: whether the textual sequence (the sequence of texts) contains the

(a) RecipeQA step text



(b) RecipeQA task title text

Figure 12: Word cloud distributions of the task title and step text in the RECIPEPLAN dataset.



(a) Wikihow step text



(b) Wikihow task title text

Figure 13: Word cloud distributions of the task title and step text in the WIKIPLAN dataset.

amount of information needed to complete the task.

- Visual-Informativeness: whether the visual sequence (the sequence of images) contains the amount of information needed to complete the task.

- Temporal Coherence: whether the multimodal sequence (the paired sequence of texts and images) meets the temporal commonsense requirements, such as a step occurring before another step instead of after.

- Plan Accuracy: whether the multimodal sequence (the paired sequence of texts and images) can successfully complete the task.

To be concrete, the annotators were asked to choose one from the two sequences by *1 - Sequence 1 is better*, *2 - Tie*, and *3 - Sequence 2 is better*. We provide the multimodal plans as follows:

**Task: How to Get Kids to Eat Healthy.**
**Sequence 1:**

Step 1: Talk to your kids about the importance of eating healthy and make sure that nutritious food is accessible and visible in the house.

Visual Plan at Step 1: [Image]

Step 2: Set rules or guidelines for what is allowed and not allowed in terms of snacks and

meals., involve children in grocery shopping and meal preparation as much as you can, lead by example by practicing healthy eating habits yourself.

Visual Plan at Step 2: [Image]

Step 8: Encourage better dietary decisions at mealtimes by sitting down together with them at the table filled with various types of wholesome foods.

Visual Plan at Step 8: [Image]

**Sequence 2:**

Step 1: Talk to your kids about the importance of eating healthy.

Visual Plan at Step 1: [Image]

Step 2: Set rules or guidelines for what is allowed and not allowed in terms of snacks and meals.

Visual Plan at Step 2: [Image]

Step 8: Reward your child for trying new foods or making healthier choices.

Visual Plan at Step 8: [Image]

We show the paired visual plan at each step below the text plan.

Please refer to our Amazon Mechanical Turk human evaluation interface for head-to-head comparison on WIKIPLAN and RECIPEPLAN in Figure 15 and Figure 14 respectively.

Figure 14: Amazon Mechanical Turk Platform. Questions Layout for Human Raters for Win-Tie-Lose Comparison on WIKIPLAN

# E  More Results

## E.1  Showcases

We show more cases in Figure 16- 19 comparing our Text-Image Prompting with powerful baselines Text-Davinci-003 + Stable-Diffusion.

## E.2  Word Cloud

In comparison with the word cloud distribution of the ground truth, we also show the word cloud of the baselines and TIP on WIKIPLAN and RECIPEPLAN.

## E.3  Failure Cases

In Figure 20, we showcase failure generation. For example, the state of the almond stays unchanged in Figure 20a, we suppose this is due to no explicit awareness of previous state change. In Figure 20b, at step 2, the generated image plan, though complemented with the text plan, loses authenticity in that the clock should not appear in a pan with carrots.

## Overview

Thank you for helping us with our research! Please read the below instruction before you go to the question.

**Instruction:**
Given the Task (e.g, Task: How to muddle), please compare two sequences of steps (Sequence 1 and Sequence 2) and **determine which sequence is better** in terms of four aspects:

- **Textual-Informativeness:** whether the textual sequence (the sequence of texts) contain the amount of information needed to complete the task.
- **Visual-Informativeness:** whether the visual sequence (the sequence of images) contain the amount of information needed to complete the task.
- **Temporal Coherence:** whether the multimodal sequence (the paired sequence of texts and images) meet the temporal commonsense requirements, such as a step occurs before another step instead of after.
- **Plan Accuracy:** whether the multimodal sequence (the paired sequence of texts and images) can successfully complete the task.
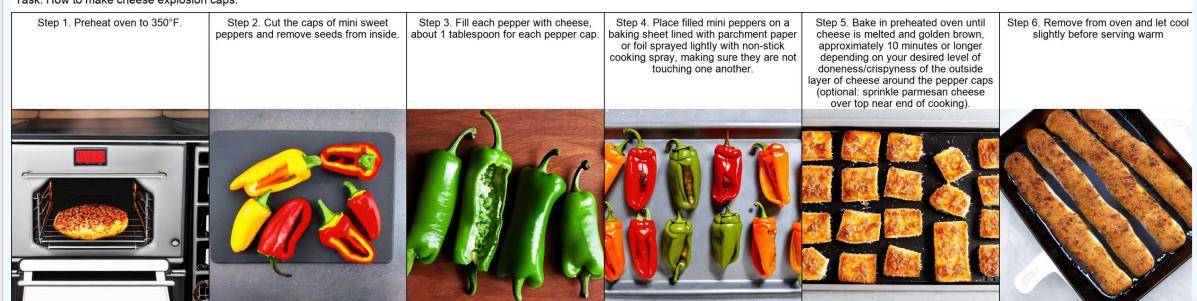
## Question

**${task_txt}**

**Sequence 1:**

Task: How to make cheese explosion caps.

| Step 1. Preheat the oven to 350°F and prepare a baking sheet lined with parchment paper or foil sprayed lightly with non-stick cooking spray. | Step 2. Cut off the tops of mini sweet peppers, removing seeds from inside each pepper cap. | Step 3. Fill each pepper cap with approximately 1 tablespoon of cheese. | Step 4. Place filled mini peppers on the prepared baking sheet, making sure they are not touching one another. | Step 5. Bake in preheated oven until cheese is melted and golden brown, approximately 10 minutes or longer depending on your desired level of doneness/crispyness of the outside layer of cheese around the pepper caps (optional: sprinkle parmesan cheese over top near end of cooking). | Step 6 Remove from oven and let cool slightly before serving warm |

**Sequence 2:**

Task: How to make cheese explosion caps.

| Step 1. Preheat oven to 350°F. | Step 2. Cut the caps of mini sweet peppers and remove seeds from inside. | Step 3. Fill each pepper with cheese, about 1 tablespoon for each pepper cap. | Step 4. Place filled mini peppers on a baking sheet lined with parchment paper or foil sprayed lightly with non-stick cooking spray, making sure they are not touching one another. | Step 5. Bake in preheated oven until cheese is melted and golden brown, approximately 10 minutes or longer depending on your desired level of doneness/crispyness of the outside layer of cheese around the pepper caps (optional: sprinkle parmesan cheese over top near end of cooking). | Step 6. Remove from oven and let cool slightly before serving warm |

**1. Select an Option: [Please only consider texts in this option.]** Which sequence contains richer Textual-Informativeness? The richer the textual sequence, the more information in the textual sequence (Step 1, Step 2, ...) is for completing the give Task .

○ 1 - Sequence 1 is better    ○ 2 - Tie    ○ 3 - Sequence 2 is better

**2. Select an Option: [Please only consider images in this option.]** Which sequence contains richer Visual-Informativeness? The richer the visual sequence, the more information in the visual sequence (the image below each Step) is for completing the given Task .

○ 1 - Sequence 1 is better    ○ 2 - Tie    ○ 3 - Sequence 2 is better

**3. Select an Option: [Please consider both texts and images in this option.]** Which sequence has more reasonable Temporal Coherence>? In other words, the reasonable multimodal (text and image) sequence obeys commonsense knowledge in real-world that certain steps should happen before others in the given Task .

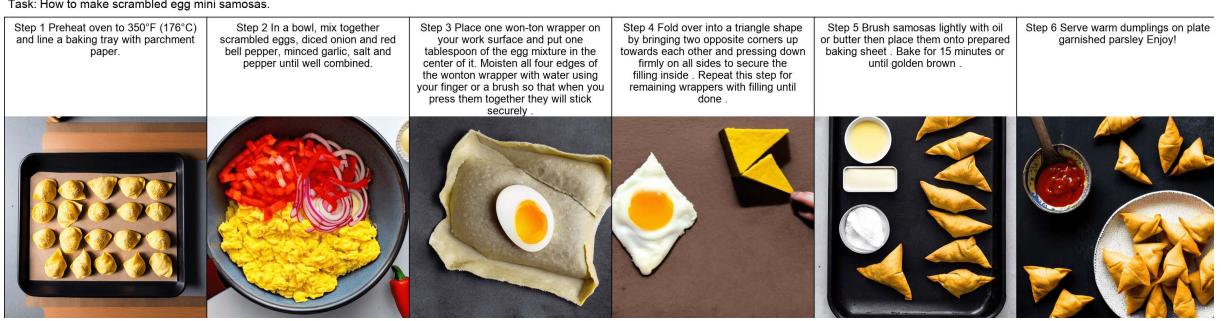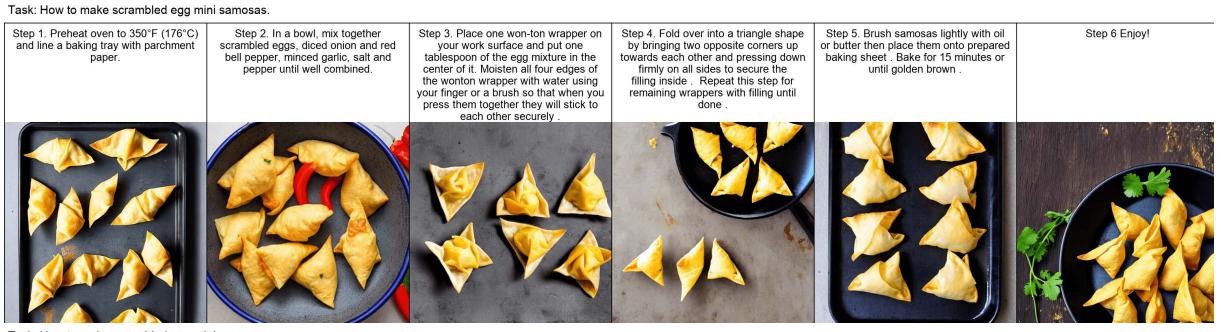○ 1 - Sequence 1 is better    ○ 2 - Tie    ○ 3 - Sequence 2 is better

**4. Select an Option: [Please consider both texts and images in this option.]** Which sequence can achieve higher Plan Accuracy>? In other words, can the multimodal (text and image) sequence more successfully complete the given Task ?

○ 1 - Sequence 1 is better    ○ 2 - Tie    ○ 3 - Sequence 2 is better

Figure 15: Amazon Mechanical Turk Platform. Questions Layout for Human Raters for Win-Tie-Lose Comparison on RECIPEPLAN.

(a) Text-Davinci-003 + Stable-Diffusion (Top) vs. TIP (Bottom)



(b) Text-Davinci-003 + Stable-Diffusion (Top) vs. TIP (Bottom)

Figure 16: More showcases of plan comparisons on RECIPEPLAN.

Task: How to make hot dogs and potatoes slachipapas with refajo colo.

| Step 1. Boil potatoes in salted water until tender | Step 2. Fry hot dogs in a skillet for about 5 minutes, stirring often | Step 3. Slice the potatoes and add to the hot dog pan with 1/4 cup of refajo colo | Step 4. Sauté all ingredients together until heated through | Step 5. Serve on warmed tortillas or buns with guacamole, sour cream and cheese |
|---|---|---|---|---|



Task: How to make hot dogs and potatoes slachipapas with refajo colo.

| Step 1. Boil potatoes in salted water until tender | Step 2. Fry hot dogs in a skillet for about 5 minutes, stirring often | Step 3. Slice the boiled potatoes and add them to the pan of frying hotdogs along with 1/4 cup of refajo colo | Step 4. Sauté all ingredients together until heated through | Step 5 Serve on warmed tortillas or buns topped guacamole, sour cream, cheese, and other desired toppings |
|---|---|---|---|---|



(a) Text-Davinci-003 + Stable-Diffusion (Top) vs. TIP (Bottom)

Task: How to make peppermint patty pudding shot.

| Step 1. Combine 1/2 cup of peppermint schnapps, 1/4 cup of Bailey's Irish cream, and 2 tablespoons of chocolate syrup in a bowl or blender | Step 2. Add 4 ounces of cold heavy whipping cream to mixture and blend until smooth | Step 3. Place 4 small glasses on a flat surface and fill with pudding shot mix | Step 4. Chill for at least 30 minutes before serving |
|---|---|---|---|



Task: How to make peppermint patty pudding shot.

| Step 1. Place three cups of chocolate pudding on a flat surface and add 1/2 cup of peppermint schnapps, 1/4 cup of Baileys Irish Cream, and 2 tablespoons of chocolate syrup | Step 2. Using a spoon or blender mix all ingredients until smooth | Step 3. Pour the mixture into four small glasses and top each with whipped cream | Step 4. Add a candy cane for garnish on each glass and chill for at least 30 minutes before serving |
|---|---|---|---|



(b) Text-Davinci-003 + Stable-Diffusion (Top) vs. TIP (Bottom)

Figure 17: More showcases of plan comparisons on RECIPEPLAN.

Task: How to Use Oregano Oil

| Step 1. Purchase oregano oil supplement capsules or pure essential oil | Step 2. Dilute the essential oil with a carrier oil such as olive, coconut, sweet almond or jojoba oils | Step 3. Test for sensitivity by applying a small amount of diluted oregano oil to the inside of your arm and wait 24 hours | Step 4. Ingest one capsule daily or mix 3-5 drops of diluted essential oil in water and drink it once per day | Step 5. Apply 2-3 drops of undiluted oregano essential oil directly on areas afflicted with infection such as cold sores twice daily | Step 6. For skin issues, massage 3-5 drops into affected area twice daily |

Task: How to Use Oregano Oil

| Step 1. Purchase oregano oil supplement capsules or pure essential oil and a carrier oil such as olive, coconut, sweet almond or jojoba oils | Step 2. Grind the oregano in a mortar and pestle to form a powder | Step 3. Dilute the essential oil with your chosen carrier oil in order to create an ointment | Step 4. Test for sensitivity by applying a small amount of diluted oregano oil ointment to the inside of your arm and wait 24 hours | Step 5. Ingest one capsule daily or mix 3-5 drops of diluted essential oil ointment in water and drink it once per day | Step 6. For skin issues, massage 3-5 drops of undiluted essential oregano oil directly on areas afflicted with infection such as cold sores twice daily |

(a) Text-Davinci-003 + Stable-Diffusion (Top) vs. TIP (Bottom)

Task: How to Make Graham Balls

| Step 1. Combine graham cracker crumbs and melted butter in a bowl. | Step 2. Mix until the ingredients are fully incorporated. | Step 3. Create small balls with the mixture and set aside on a plate or baking sheet. | Step 4. Place condensed milk, peanut butter, and chocolate chips into a separate bowl, mix together until combined evenly throughout the ingredients.. | Step 5. Dip each of the graham ball mixtures into the condensed milk mixture then roll them around to cover completely in coating before placing back onto plate or baking sheet once more.. | Step 6. Freeze for 1-2 hours to allow balls to harden completely before serving/eating! |

Task: How to Make Graham Balls

| Step 1. Combine graham cracker crumbs, melted butter and spices in a bowl and mix until fully incorporated. | Step 2. Crack eggs into the mixture, mixing everything together evenly. | Step 3. Create small balls with the mixture and set aside on a plate or baking sheet. | Step 4. In a separate bowl combine condensed milk, peanut butter, chocolate chips and peanuts together until blended evenly throughout the ingredients.. | Step 5 Dip each of the graham ball mixtures into the condensed milk-peanut butter-chocolate chip-peanuts mixture then roll them around to cover completely in coating before placing back onto plate or baking sheet once more.. | Step 6 Slice oranges into wedges and place one wedge on top of each ball before freezing for 1-2 hours allowing them harden completely before serving/eating! |

(b) Text-Davinci-003 + Stable-Diffusion (Top) vs. TIP (Bottom)

Figure 18: More showcases of plan comparisons on WIKIPLAN.

Task: How to Make Pizza Lunchables

| Step 1. Preheat oven to 400°F. | Step 2. Unwrap two Lunchables pepperoni pizzas and place them on a baking sheet lined with parchment paper or aluminum foil. | Step 3. Bake for 8 minutes, then remove from the oven and let cool for 5 minutes before serving. | Step 4. Cut each pizza into four slices using a pizza cutter or sharp knife, then serve with your favorite dipping sauce such as ranch dressing, barbecue sauce, ketchup, etc.. |
|---|---|---|---|



Task: How to Make Pizza Lunchables

| Step 1. Unwrap two Lunchables pepperoni pizzas and place them on a baking sheet lined with parchment paper or aluminum foil. | Step 2. Preheat oven to 400°F and bake for 8 minutes. | Step 3. Remove from the oven, let cool for 5 minutes before cutting each pizza into four slices using a pizza cutter or sharp knife. | Step 4. Place the slices in lunch box with chips and milk, then serve with your favorite dipping sauce such as ranch dressing, barbecue sauce, ketchup etc.. |
|---|---|---|---|



(a) Text-Davinci-003 + Stable-Diffusion (Top) vs. TIP (Bottom)

Task: How to Make Ice Balls

| Step 1. Fill an ice cube tray with water | Step 2. Place the tray in the freezer and freeze until solid | Step 3. Remove the frozen cubes from the tray | Step 4. Arrange three cubes together in a round shape and press them together lightly to form a ball shape | Step 5. Dip the ball into warm water for a few seconds to make it easier to handle | Step 6. Roll each ball on a flat surface or between two spoons until completely smooth and round |
|---|---|---|---|---|---|



Task: How to Make Ice Balls

| Step 1. Fill an ice cube tray with water and place it in the freezer until solid. | Step 2. Take out the frozen cubes from the tray and arrange them into a round shape, pressing lightly to form a ball shape. | Step 3. Dip the ball into warm water for a few seconds to make it easier to handle. | Step 4. Roll each ball on a flat surface or between two spoons until completely smooth and round, like rolling sugar cubes on a countertop or chocolate balls on a metal tray. | Step 5. a bowl of warm water, a person dipping an ice ball into the water, and an ice ball. | Step 6. a smooth, flat surface, like a cutting board, or two spoons, and a few ice balls being rolled between them. |
|---|---|---|---|---|---|



(b) Text-Davinci-003 + Stable-Diffusion (Top) vs. TIP (Bottom)

Figure 19: More showcases of plan comparisons on WIKIPLAN.

Task: How to make homemade raw almond butter.

| Step 1. Place almonds in a bowl and soak overnight | Step 2. Preheat oven to 350°F and spread soaked almonds on a baking sheet | Step 3. Roast for 10-15 minutes, stirring occasionally until lightly browned | Step 4. Remove from oven and let cool completely | Step 5 . Place cooled roasted almonds in food processor or high-speed blender; blend until creamy, stopping to scrape down sides as needed | Step 6 Add salt (optional) and continue blending until desired consistency is reached |
|---|---|---|---|---|---|



(a) Failure multimodal plans generated by our Text-Image Prompting (TIP).

Task: How to make carrot and swede potch mash.

| Step 1. Peel and chop carrots and swedes into small cubes | Step 2. Boil the vegetables in a saucepan with boiling salted water for 15-20 minutes until tender | Step 3. Place the cut up carrots into a wooden bowl before adding them to the pan of boiling water along with the swede cubes | Step 4. Drain off excess liquid from the cooked vegetables | Step 5. Mash the cooked vegetables with butter, salt and pepper to taste |
|---|---|---|---|---|



(b) Failure multimodal plans generated by our Text-Image Prompting (TIP).

Figure 20: We showcase failure cases of our Text-Image Prompting on generating multimodal plans on both datasets.