

MATEO: A Multimodal Benchmark for Temporal Reasoning and Planning in LVLMs

Anonymous ACL submission

Abstract

AI agents need to plan to achieve complex goals that involve orchestrating perception, sub-goal decomposition, and execution. These plans consist of ordered steps structured according to a Temporal Execution Order (TEO, a directed acyclic graph that ensures each step executes only after its preconditions are satisfied. Existing research on foundational models' understanding of temporal execution is limited to automatically derived annotations, approximations of the TEO as a linear chain, or text-only inputs. To address this gap, we introduce MATEO (Multimodal Temporal Execution Order), a benchmark designed to assess and improve the temporal reasoning abilities of Large Vision Language Models (LVLMs) required for real-world planning. We acquire a high-quality professional multimodal recipe corpus, authored through a standardized editorial process that decomposes instructions into discrete steps, each paired with corresponding images. We collect TEO annotations as graphs by designing and using a scalable crowdsourcing pipeline. Using MATEO, we evaluate six state-of-the-art LVLMs across model scales, varying language context, multimodal input structure, and fine-tuning strategies.

1 Introduction

The era of autonomous agents with natural language interfaces has recently gained attention from both academia and industry. AI agents have been successfully adopted across many fields, including finance (Fatemi and Hu, 2024; Zhang et al., 2024), software development (Suri et al., 2023; Mo et al., 2025), and networking (Zhang and Zhu, 2023). Agents are autonomous systems that can accomplish goals with little or no human supervision, a concept that dates back to the earliest days of artificial intelligence (Weizenbaum, 1966; Wooldridge and Jennings, 1995). Several frameworks have been proposed to describe the architec-

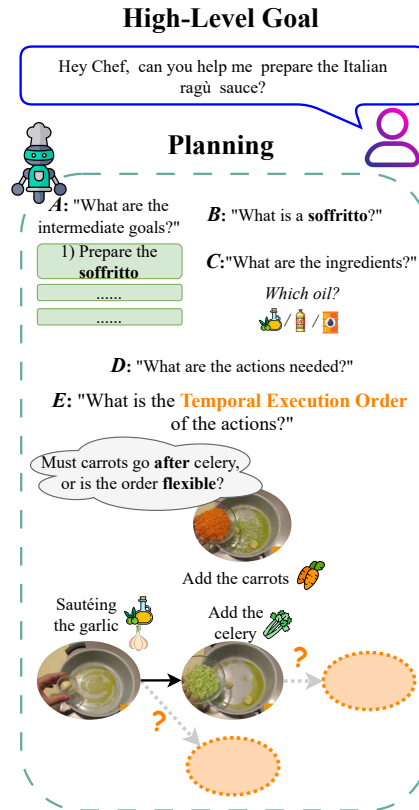


Figure 1: Example of an AI agent’s planning process and inherent uncertainties for a natural-language goal. Questions A–D decompose the high-level goal into executable actions, while E infers their Temporal Execution Order (TEO) as a directed acyclic graph. Each step introduces uncertainty, producing multiple possible paths, some correct, others wrong.

ture of these systems, including the BDI (belief-desire-intention) model (Rao et al., 1995) or, more recently, the perceive-reason-act-learn paradigm in agentic AI (Raheem and Hossain, 2025). Despite differences in terminologies and definitions, these frameworks generally decompose the agent behavior into environmental sensing, reasoning and understanding, planning, and action execution.

Planning is one of the most complex components, as it requires mapping high-level goals to

043
044
045
046
047
048
049
050
051
052

053 executable action sequences. Its complexity in- 105
054 creases when the input is natural language pro- 106
055 cessing alongside other sources, such as images or 107
056 sensor signals. Classical symbolic planners, such 108
057 as STRIPS (Fikes and Nilsson, 1971), represent 109
058 as logic conditions to be satisfied through actions 110
059 defined by explicit pre- and post-conditions. Al- 111
060 though these methods offer strong guarantees of 112
061 soundness and completeness and are explainable by 113
062 design, they depend on hand-crafted domain mod- 114
063 els and thus scale poorly to dynamic environments 115
064 (Acharya et al., 2025). 116

065 Foundation language models offer a promising 117
066 alternative, as they can directly manipulate goals 118
067 and actions expressed in natural language without 119
068 requiring a formal symbolic description, making 120
069 them flexible to adapt to dynamic environments 121
070 (Raheem and Hossain, 2025). To generate an ef- 122
071 fective plan, a language model first needs to parse 123
072 the goal and available actions, while addressing the 124
073 inherent ambiguity of natural language and relative 125
074 uncertainties, by following a process similar to that 126
075 depicted in Figure 1. Specifically, it then needs to 127
076 reason over this inferred information, together with 128
077 parametric knowledge or retrieved external data, to 129
078 construct a sequence of actions that achieves the 130
079 desired goals. 131

080 Current research on Large Language Mod- 132
081 els (LLMs) and Large Vision-Language Models 133
082 (LVLMs) as planners has shown serious limitations 134
083 when dealing with complex tasks such as travel, 135
084 flight, and calendar planning (Xie et al., 2024; Ji 136
085 et al., 2025; Zheng et al., 2024). Nonetheless, most 137
086 assessments of these focus on plan-level metrics, 138
087 such as the success rate, providing little insight 139
088 into the fine-grained sources of errors. While plan- 140
089 ning is generally considered a PSPACE-complete 141
090 problem (Bylander, 1994), understanding whether 142
091 these models can reliably infer the preconditions 143
092 and effects of actions, which is a fundamental abil- 144
093 ity for any planning system, might help improve 145
094 their performance and deepen our understanding of 146
095 their actual reasoning abilities. 147

096 A possible way to investigate these abilities is to 148
097 evaluate the models’ performance on the Temporal 149
098 Execution Order (TEO) task (Chambers, 2013; 150
099 Derczynski, 2017; Vashishtha et al., 2020). In this, 151
100 given two steps representing two actions, the model 152
101 has to infer their temporal relation, such as which 153
102 action should occur first or second, or if the actions 154
103 are order-independent. Successful resolution of 155
104 this requires reasoning over the actions’ implicit

105 pre- and post-conditions. Procedural texts, such as 106
107 recipes or instruction books, provide an ecologi- 107
108 cally valid and scalable source of such instances, 108
109 as they explicitly encode action sequences as plans 109
110 expressed through stepwise natural-language in- 110
111 structions. However, existing benchmarks rely on 111
112 text-only procedural text (Mori et al., 2014; Lal 112
113 et al., 2024), automatically derived annotations, or 113
114 approximate TEO as a strictly linear chain (Wu 114
115 et al., 2022; Lu et al., 2024; Qiu et al., 2025). This 115
116 limits the interpretation of the model’s actual per- 116
117 formance on the TEO and the insights into whether 117
118 LVLMs can function as *world models*, i.e., systems 118
119 capable of grounding plans on world observations 119
(Chen et al., 2025). 120

120 In this work, we introduce MATEO 120
121 (Multimodal Temporal Execution Order), a 121
122 publicly available benchmark designed for evalu- 122
123 ating and improving the multimodal temporal 123
124 reasoning abilities of LVLMs. MATEO assesses 124
125 a model’s capability to determine the Temporal 125
126 Execution Order (TEO) of a sequence of multi- 126
127 modal actions, a fundamental building block for 127
128 planning in real-world contexts. The benchmark 128
129 comprises 300 high-quality professional recipes 129
130 from a well-known Italian recipe website¹. Each 130
131 recipe is organised as a sequence of steps, with 131
132 every step containing both a textual description 132
133 and an image illustrating the action or its outcome, 133
134 ensuring quasi-perfect semantic alignment between 134
135 text and visual content. Different from most 135
136 existing human-annotated recipe corpora (Mori 136
137 et al., 2014; Yamakata et al., 2020; Pan et al., 137
138 2020), which rely on internal annotators for TEO 138
139 labelling, we employ crowdsourcing, providing a 139
140 scalable and reproducible annotation methodology. 140
141 The annotation process has produced a Directed 141
142 Acyclic Graph (DAG) for each recipe, capturing 142
143 the pre- and post-condition dependencies (*edges*) 143
144 of each action (*nodes*). Using this benchmark, we 144
145 evaluate six open and closed-sourced LVLMs with 145
146 various prompting strategies, including in-context 146
147 learning (Brown et al., 2020), chain of thought 147
148 (Wei et al., 2022), and self-reflection (Madaan 148
149 et al., 2023), and fine-tuning approaches. Our 149
150 findings reveal that while some models achieve 150
151 state-of-the-art results with multimodal inputs, 151
152 the majority still struggle to effectively leverage 152
153 both modalities. Furthermore, even the best 153
154 performing models only achieve 0.69 accuracy, 154

¹<https://www.giallozafferano.it/>

155 highlighting suboptimal capabilities in the TEO
156 task and motivating future adoption of MATEO
157 for developing novel methods to enhance temporal
158 reasoning and real-world planning.

159 In summary, the main contributions of this paper
160 are:

- 161 • Release MATEO, a publicly available² multi-
162 modal benchmark for evaluating and advancing
163 the temporal reasoning abilities of LVLMs
164 through the Temporal Execution Order (TEO)
165 task;
- 166 • A systematic evaluation of six open- and
167 closed-source LVLMs using in-context learning
168 (ICL), chain-of-thought prompting (CoT),
169 self-reflection, and fine-tuning approaches
170 while varying the input modalities;
- 171 • A set of novel guidelines for extending the
172 benchmark to new recipes, other languages,
173 and domains, ensuring scalability and repro-
174 ducibility.

175 2 Literature Review

176 **Planning** The task of planning is to find a sequence
177 of actions that, when executed, transitions a sys-
178 tem from an initial state to a desired goal state
179 (Valmeekam et al., 2023). Recent work has in-
180 vestigated LLMs’ planning abilities by evaluat-
181 ing them on classical planning problems, such
182 as Blocks World (Gupta et al., 1991), showing
183 their limitations in symbolic planning (Valmeekam
184 et al., 2023; Kambhampati et al., 2024). However,
185 these planning problems might be overspecific and
186 formalized, impeding the model from leveraging
187 any parametric knowledge, and are described with
188 hand-crafted actions and states, impacting the flex-
189 ibility to adapt to dynamic environments. In this
190 regard, several benchmarks in which initial and
191 goal states are written in natural language have
192 been proposed. They cover various domains such
193 as the planning of trips, meetings, and calendars
194 (Zheng et al., 2024; Xie et al., 2024). Nonetheless,
195 these are only for sequential planning, i.e., without
196 accounting for possible execution of simultaneous
197 actions, and most of the proposed assessments are
198 based on plan-level metrics. A few studies repre-
199 sent the plan as a graph, but they rely on machine-
200 generated annotations (Lin et al., 2024). Although

²We have received clearance from GialloZafferano pub-
lisher to share the data and its annotation for research pur-
poses.

201 multimodal planning is a relatively new topic, a few
202 benchmarks for investigating the planning abilities
203 of LVLMs have been proposed, such as MPCC (Ji
204 et al., 2025), EgoPlan-Bench2 (Qiu et al., 2025),
205 ALFWorld (Li et al., 2024b), and procedural text
206 corpora Wu et al. (2022); Lu et al. (2024). How-
207 ever, all of these are limited to a linear chain of
208 action due to goals with low uncertainty or a lack
209 of an ad-hoc annotation.

210 **TEO in Recipes** Recipes constitute a chal-
211 lenging multimodal planning domain and a rich
212 procedural-text genre for the TEO task, as their
213 steps often form parallelizable sub-plans with par-
214 tial ordering constraints. Among the earliest efforts
215 to model this planning complexity are the Japanese
216 Recipe Flow Graph (Mori et al., 2014) and the En-
217 glish Recipe Flow Graph (Yamakata et al., 2020)
218 corpora, which model directly the dependencies
219 and effects of each action expressed by verbs, simi-
220 larly to a dependency parser. Step-level dependen-
221 cies have been extracted from the English Recipe
222 Flow Graph to create the CaT-Bench corpus (Lal
223 et al., 2024). Although a multimodal recipe corpus
224 following a similar logic has been proposed (Pan
225 et al., 2020), no public multimodal DAG corpus
226 is currently available. Existing resources rely on
227 non-curated web/video data, leading to inconsis-
228 tent formats³, and their annotations are produced
229 by non-crowdworkers, limiting scalability and re-
230 producibility.

231 3 Methodology

232 To evaluate the reasoning abilities required for plan-
233 ning in LLMs and LVLMs, we introduce the Tem-
234 poral Execution Order (TEO) task. We formulate
235 TEO as a three-way classification problem over
236 the temporal execution relations *before* and *after*,
237 and *independent*. These relations are semantically
238 aligned with the *before*, *after*, and *simultaneous*
239 temporal relations defined in Allen’s interval al-
240 gebra (Allen, 1983). However, unlike Allen’s al-
241 gebraic relations, which are derived from compar-
242 isons of interval start and end points, our formula-
243 tion infers temporal execution order from logical
244 dependencies rather than explicit time boundaries.
245 In particular, if action B requires a condition that
246 is produced by the post-condition (aka effect) of
247 action A, then A must be executed *before* B (equiv-
248 alently, B must occur *after* A). Conversely, if A and
249 B do not depend on each other’s post-conditions,

³e.g., variable image counts per step, mixed unit systems.

including indirect or transitive dependencies (e.g., $A \rightarrow C \rightarrow B$, where B requires C and C requires A), then the two actions are considered independent, meaning they can be executed in any order. Given a set of action instances \mathcal{A} and their dependency relations \mathcal{E} , TEO can be represented as the graph $G = (\mathcal{A}, \mathcal{E})$, where \mathcal{A} are the nodes, and \mathcal{E} are the edges. Specifically, this graph is a directed acyclic graph (DAG) because the dependencies have an orientation and any cycle would generate a paradox, making the plan infeasible (e.g., an action dependent on its own output).

Similar to Roccabruna et al. (2024) and Lal et al. (2024), for each class, we pose a corresponding ternary (Yes/No/I don’t know) question, as presented in Table 1.

TEO class	Question
Before	Must Step A be executed before Step B?
After	Must Step A be executed after Step B?
Independent	Can Step A and Step B be executed in parallel ?

Table 1: Temporal Execution Order (TEO) classes and their corresponding ternary questions.

To investigate whether the models can answer these questions, we formalise the following input sequence. Since our corpus consists of recipes, we refer to actions as steps, following the natural structure of a recipe. In our corpus, a step is composed of an image I and a textual part T , formally $S_j = I_j \oplus T_j$. Thus, given two steps, A and B, the model is tasked to answer the three TEO questions of Table 1, denoted as Q_i , using the following input sequence:

$$C \oplus S_A \oplus S_B \oplus Q_1 \oplus Q_2 \oplus Q_3$$

where C is the grounding context, \oplus denotes concatenation with a newline (“\n”). Marker tokens (“Step A picture:”, “Step A description:”, etc.) precede each image and text component to further guide the model.

The TEO class is then retrieved by parsing the model’s output by assigning class $i \in \{before, after, independent\}$, if the model answers “Yes” to the corresponding Q_i TEO question and “No” to both other Q_j for $j \neq i$. All the other combinations, such as answering “Yes” to multiple questions or responding “I don’t know”, are

assigned to the class *Other*. This captures inconsistent behavior, such as simultaneously affirming the *before* and *after* questions, as well as uncertainty.

In Chain of Thought experiments (Wei et al., 2022), the model is prompted to directly predict one of the three TEO classes. The reason for this was to simplify and shorten the demonstrations provided in the prompt.

4 The Benchmark

Our requirements for building MATEO were to identify a form of procedural text in which the textual and visual components are maximally semantically aligned, enabling reliable multimodal reasoning over execution order. We have found that these criteria are met by the professionally edited recipes on GialloZafferano, a well-known Italian-language recipe website. These recipes offer several advantages: i) they follow an editorially-curated narrative and structural format; ii) the images are professionally produced, visually consistent, and captured from a near-uniform top-down perspective, minimizing viewpoint variability; and iii) each recipe is explicitly segmented into an ordered sequence of steps. Each step typically describes a single executable action, and the accompanying image depicts either the action’s outcome or the action in progress. Together, these properties reduce multimodal noise sources and enable cleaner corpus utilization ⁴.

To build MATEO, we employed human annotators recruited through Prolific⁵, a crowdsourcing platform. To balance crowdworker cognitive load, we divided the annotation process into batches of five recipes, with an average completion time of approximately 50 minutes per batch. We considered only native Italian speakers, given that the recipes are in Italian. We set the compensation to £9.60 per hour, which falls within the platform’s recommended range.

The annotators were provided with a list of steps in the order of the original recipe. They were tasked with linking these recipe steps into a directed acyclic graph by following the TEO rela-

⁴Because GialloZafferano’s recipes are copyright-protected and not publicly licensed for AI benchmarking, the publisher agreed to support our research initiative by providing a subset of approximately 300 recipes for release to the research community. The recipes were randomly selected via stratified sampling over dish category (e.g., appetizers, main courses, desserts) to ensure balanced coverage across recipe types.

⁵<https://www.prolific.com/>

tions by an ad-hoc UI (Figure 8 and Figure 7 in the Appendix, respectively). To explain the task concisely and effectively, we provided guidelines, a short demo video, and multimodal instructions (textual and visual), along with explanatory examples. The annotation platform and the guidelines⁶ are shown in Appendix A.2.

Additionally, to ensure data quality, annotators were required to pass a qualification test before starting the actual task. Only those who reached at least 65% accuracy on two test recipes were allowed to continue with the actual annotation work. The accuracy is computed as if it were a multi-label classification task, where the examples are the nodes and the labels are the outgoing nodes⁷. The threshold corresponds to the 75th percentile of the score distribution computed from testing 25 crowdworkers who were not involved in the actual annotation.

In total, 54 annotators⁸ annotated 315 recipes, collected in 63 batches. We computed the annotation agreement using MASI distance metrics (Pas-sonneau, 2006), which are specifically designed for measuring agreement on multi-labelling classification tasks. The annotators reached an agreement score of 0.85 (substantial agreement)⁹. Moreover, we have manually checked the annotations for specific error patterns, such as graphs with isolated steps or graphs that were fully linear. At the end of this process, we discarded 15 recipes due to poor annotation quality.

The resulting corpus contains 300 Italian recipes. Preliminary experiments show that some LVLMS perform worse on this task in Italian than in English, probably because these models are trained primarily on English (Bai et al., 2025). For this, we translate the recipe steps from Italian to English using Llama-3-8B (Grattafiori et al., 2024), given its relatively small model size and acceptable performance. Before translating the full corpus, we manually verified a subset of the translations and

⁶They will also be shared on GitHub along with the code used to conduct our experiments.

⁷To give a better visualization the data structure is $n_1 : [n_2, n_3], \dots, n_2 : [\dots]$, where n_2 and n_3 are the outgoing nodes and the “labels” of n_1 .

⁸Some of the annotators participated multiple times. This is because we experienced a shortage of annotators willing to perform our task, probably due to a high rejection rate.

⁹The agreement was computed on the examples used for qualifying the annotators. Although this may give an overestimation of the actual agreement, it was the best trade-off as we could not afford to have an overlap on the annotation due to the high amount of time needed to complete a batch.

	Train	Valid	Test
# Recipes	210	30	60
# Steps	3273	477	936
AVG Steps	15.6 ± 3.9	15.9 ± 3.7	15.6 ± 3.8
Branching Factor	1.12 ± 0.43	1.16 ± 0.42	1.11 ± 0.38

Table 2: Splits of the annotated corpus used for training and evaluating LVLMS’ capabilities. The table reports the number of samples, the average number of steps, and the branching factor.

found them to be of high quality. All the following experiments are conducted on the English translations. The dataset is split into training, validation, and test sets using stratified sampling by number of steps, with a 70%/10%/20% split. The statistics of the splits are reported in Table 2.

5 Experimental settings

5.1 Models

We select six LVLMS that vary in model size, with selection primarily based on their support for sequences of images in the input. Our evaluation includes both open-source, namely, Qwen2.5-VL-7B, Qwen2.5-VL-72B (Bai et al., 2025), LLaVA-OneVision-7B (Li et al., 2024a), InternVL3.5-8B, InternVL3.5-38B (Wang et al., 2025), and the closed-source model GPT-5.1¹⁰. We could evaluate GPT-5.1 on only 20% of the test set, selected via stratified sampling, due to budget constraints.

5.2 Dataset

Given the annotated DAG for each recipe, we infer the TEO classes. In a real-world setting, determining temporal relations would require comparing all pairs of steps, resulting in n^2 comparisons (876,096 pairs in our test set). To make evaluation tractable, we restrict comparisons only to annotated dependencies: step A must be executed *before* step B if there is a directed edge from A to B. The class *after* is the inverse of *before*, therefore, we have obtained the ground truth examples by swapping the steps A and B. Steps A and B are *independent* if there is no directed path from A to B or from B to A. The number of resulting relations is 1573 for the *independent* class, 993 for *before*, and 933 for the *after* classes.

We evaluate all models on two versions of the dataset. The first preserves the original annotation

¹⁰<https://platform.openai.com/docs/models/gpt-5.1>

order, yielding *before* and *independent* labels. The second reverses each step pair, mapping *before* to *after*, while *independent* labels remain unchanged. With this, we can assess whether the model actually understands the relations among the steps, or whether it relies on biases learnt during pre-training.

5.3 Prompt settings

We evaluate the planning abilities of LVLMs on the TEO task under different modalities and grounding contexts. Specifically, we compare text-only, image-only, and multimodal inputs, as described in Section 3. Below, we describe the grounding contexts we used to evaluate the models. The example for each prompt template can be found in Table 5 in the Appendix A.2.

Baseline The model is only conditioned on two recipe steps and the three TEO questions.

Instructions We provide the model with a set of rules that explain the meaning of the TEO classes and instruct the model to respond to each question. The rules are general for the TEO task and are not specific to the cooking domain.

In-Context Learning (ICL) The model is conditioned on the instructions described previously and three in-context examples for each TEO class. Each example consists of a pair of steps, the TEO questions, and the correct answers, accompanied by brief explanations to encourage the model to emulate the latent reasoning process. To limit context length and computation time, all in-context examples are text-only. We prepend a modality-specific instruction noting that examples contain only text descriptions, while the model input may also include images.

Chain-of-Thought (CoT) Correctly predicting the TEO class requires understanding whether the outcome of one step is a prerequisite for completing the other. In this experiment, we explicitly ask the model to identify pre- and post-conditions of each step and then reason about their dependency. After this reasoning process, the model is tasked to directly predict the TEO class.

Self-Reflection At the end of the CoT, the model is instructed to reflect on its answer by reviewing any inconsistencies or errors in the reasoning process, and then give the final answer.

5.4 Fine-Tuning

Based on the results of the prompting experiments, we have selected the best-performing small-scale

model, InternVL3.5-8B, for fine-tuning. We fine-tune the model with LoRA (Hu et al., 2021) on image-only, text-only, and image-text modalities following the *Instructions* prompt task. To prevent label imbalance, we swap only dependent step pairs to generate examples for the *after* class. For independent step pairs, swapping does not change the label and would therefore introduce duplicate examples, potentially biasing the model toward the *independent* class. This resulted in 3,499 examples for each of the *before* and *after* classes, and 4,969 for *independent*. Examples of model inputs are provided in Table 6 in the Appendix A.1 along with all the hyperparameters.

6 Results and Error Analysis

To evaluate models’ performance in predicting temporal execution order, we consider both the original and the reverse order of each step pair. This setting reflects realistic scenarios in which the true execution order may be unknown. Accuracy is reported based on a consistency criterion: a prediction is considered correct only if the model remains consistent across both the original and swapped step orders. Specifically, for dependent pairs, the model must predict *before* in the original order and *after* when swapped. For independent pairs, the model must predict *independent* in both configurations.

Table 3 reports accuracy scores across prompting strategies, and after fine-tuning. All results are obtained from a single evaluation using greedy decoding to ensure reproducibility. Starting with the prompting strategies, we observe that in the baseline setting, all models perform poorly, with near-zero accuracy, indicating limited zero-shot ability on the TEO task. Adding task instructions, in-context learning, or chain-of-thought improves performance for most models, but the effectiveness of each setting depends on the model; thus, there is no universal best strategy. Furthermore, results show that self-reflection sometimes leads to lower performance compared to CoT, suggesting a limitation of this prompting strategy for some models.

Regarding individual model performance, Qwen2.5-VL-72B achieves the highest overall accuracy among the open-source LVLMs. When using ICL with image and text inputs, it attains an accuracy of 0.68. Under the same prompting strategy, InternVL3.5-38B achieves its best performance using text-only inputs, reaching 0.59 accuracy. Among models with fewer than

Model	Modality	Baseline	Instructions	ICL	CoT	Self-Reflect.
Qwen2.5-VL-72b	Image-only	0.08	0.55	0.58	0.38	0.34
	Text-only	0.14	0.53	0.57	0.50	0.50
	Image+Text	0.17	0.61	0.68	0.59	0.54
InternVL3.5-38b	Image-only	0.01	0.51	0.53	0.49	0.44
	Text-only	0.10	0.50	0.59	0.37	0.39
	Image+Text	0.09	0.51	0.49	0.50	0.47
InternVL3.5-8b	Image-only	0.15	0.58	0.39	0.14	0.19
	Text-only	0.10	0.31	0.46	0.31	0.32
	Image+Text	0.10	0.34	0.32	0.33	0.24
Qwen2.5-VL-7b	Image-only	0.01	0.05	0.04	0.09	0.05
	Text-only	0.01	0.11	0.08	0.15	0.09
	Image+Text	0.02	0.10	0.03	0.17	0.21
LLaVA-OneVision-7b	Image-only	0.00	0.00	0.00	0.02	0.00
	Text-only	0.00	0.00	0.07	0.22	0.09
	Image+Text	0.00	0.03	0.04	0.09	0.00
GPT-5.1	Image-only	0.13*	0.54*	0.63*	0.71*	0.71*
	Text-only	0.07*	0.27*	0.27*	0.62*	0.62*
	Image+Text	0.29*	0.58*	0.7*	0.73*	0.74*
InternVL3.5-8B _{Fine-tuned}	Image-only	-	0.68	-	-	-
	Text-only	-	0.61	-	-	-
	Image+Text	-	0.69	-	-	-

Table 3: Performance of LVLMs on the TEO task under different prompt strategies and fine-tuning results. We report accuracy by considering the consistency of predictions across the original and swapped orders of steps. The results in bold indicate the best results per model. *Results achieved on 20% of the test set only.

10B parameters, Qwen2.5-VL-7B and LLaVA-OneVision-7B substantially underperform, with accuracy at or near zero in most settings. In contrast, InternVL3.5-8B achieves performance comparable to larger-scale models, reaching 0.58 accuracy when prompted with task instructions and image-only inputs. While the GPT-5.1 results are not directly comparable¹¹ (we evaluate it only on 20% of the test set), our findings suggest that it may perform better than the open-source models, achieving 0.74 accuracy with self-reflection and multimodal inputs. In contrast, some models fail to self-reflect or produce a final answer, with failure rates ranging from 2.71% for Qwen2.5-VL-72B to 81.6% for LLaVA-OneVision-7B, compared to just 0.03% for GPT-5.1. When self-reflection is successful, models largely preserve their original predictions (98.5% for GPT-5.1 and 97.0% for Qwen2.5-VL-72B). Notably, GPT-5.1 makes only correct self-corrections (0.9%), whereas Qwen2.5-VL-72B makes nearly equal numbers of

¹¹Not comparable but indicative as Qwen2.5-VL-72B on the same partition with ICL Image+Text achieves 0.66 compared to 0.68 on the whole test set.

correct and incorrect self-corrections (0.14% and 0.13%), suggesting a possible reason for GPT-5.1’s stronger gains under this strategy.

Although multimodal input leads Qwen2.5-VL-72b and GPT-5.1 to achieve state-of-the-art results, not all models show the same ability to leverage this modality. Indeed, other models struggle to incorporate multimodal information, and we observe a drop in performance compared to single modalities. A possible explanation could be that the scale of the language decoder, with small-scale models, fails to effectively attend to different modalities.

Regarding the fine-tuning experiments, InternVL3.5-8B achieves 0.69 of accuracy, which is the highest among other experiments with open-source models. When restricted to a single modality, InternVL3.5-8B also outperforms all other open-source models in that modality. Notably, image-only fine-tuning yields better performance than text-only fine-tuning. Some possible explanations are that the model cannot fully leverage textual input alone, or that images more precisely capture the underlying actions, whereas textual descriptions may contain superflu-

Model	Class	Baseline	Instructions	ICL	CoT	Self-Reflect.
GPT-5.1	Before	0.8*	0.82*	0.79*	0.78*	0.79*
	Independent	0.39*	0.79*	0.85*	0.87*	0.87*
	After	0.65*	0.76*	0.7*	0.69*	0.72*
Qwen2.5-VL-72b	Before	0.6	0.76	0.79	0.79	0.78
	Independent	0.21	0.79	0.85	0.86	0.86
	After	0.5	0.71	0.71	0.42	0.43
InternVL3.5-8b	Before	0.6	0.35	0.63	0.71	0.67
	Independent	0.0	0.73	0.63	0.74	0.69
	After	0.33	0.51	0.47	0.18	0.16

Table 4: Per-class F1 scores for the TEO task. Scores above the dashed line are computed on the original step order, and those below on the swapped order. Models in zero-shot show bias toward *before* and *after* relations, lowering *independent* performance. Prompting improves *independent* and *after* scores; however, *after* remains challenging, suggesting difficulties in reasoning over reversed temporal order. *Results achieved on 20% of the test set only.

ous information, such as references to other steps, or lack sufficient contextual information.

We analyze performance across individual TEO classes to identify the most challenging cases. Table 4 reports F1 scores for the *before*, *independent*, and *after* classes. Because the ground-truth labels for the *after* class are defined on the swapped step order, we compute F1 scores for this class on the swapped dataset. In contrast, scores for the *before* and *independent* classes are computed using the original step order. For clarity, we present the selection of three models, all under multimodal inputs; similar trends are observed across all models. Under the baseline zero-shot setting, models exhibit a bias toward *before* and *after* relations, lowering performance on the *independent* class. Performance on the *independent* class improves substantially with task instructions, ICL, and CoT prompting, while gains for the *before* class are more moderate. For Qwen2.5-VL-72B, CoT yields relative improvements of 32% and 310% for the *before* and *independent* classes with respect to the baseline, respectively. Moreover, the models struggle in predicting the *after* class, particularly under the baseline and CoT settings. This suggests a bias inherited from pre-training, where models have difficulty reasoning under reversed TEO. When conditioned on task instructions, Qwen2.5-VL-72B achieves an F1 score of 0.71 on the *after* class, while InternVL3.5-8B reaches a maximum F1 of 0.51.

Higher per-class scores compared to the consistency-based accuracy indicate that models struggle to produce consistent predictions across the original and swapped step orders.

Larger-scale models, namely Qwen2.5-VL-72B and InternVL3.5-38B, exhibit greater consistency in their predictions. In particular, they produce consistent responses in 75% and 68% of cases for the *independent* class, and in 69% and 73% of cases for *before–after* relations, respectively. In contrast, smaller-scale models struggle substantially with consistency, particularly for *before–after* relations; for example, InternVL3.5-8B achieves only 34% consistency under its best-performing prompting setting. This further highlights that models struggle to reason under the reversed temporal order. Fine-tuning improves prediction consistency, with InternVL3.5-8B achieving 76% consistency for *before–after* pairs and 73% for *independent* pairs.

6.1 Conclusion

We introduced MATEO, a publicly available multimodal benchmark designed to evaluate and advance the multimodal temporal reasoning and thus planning abilities of LVLMs via the Temporal Execution Order task. We provided comprehensive guidelines for extending the benchmark to new recipes, languages, and domains, enabling scalable, reproducible use for both evaluation and model development. We evaluated six open- and closed-source models and demonstrated that multimodality is critical for successfully performing the TEO task. Nevertheless, even the best-performing model achieved only modest performance, highlighting a possible explanation for the limited planning capabilities of current LVLMs. We believe that MATEO can serve as a foundation for future research and improvements in multimodal temporal reasoning and real-world planning.

623 Limitations

624 The main limitation of this work lies in the prompt-
625 ing strategies. We observed that models often
626 interpreted instructions differently, requiring mi-
627 nor model-specific adjustments, such as rephrasing
628 the output format or clarifying the reasoning in-
629 structions, and preventing a fully uniform, directly
630 comparable evaluation. Specifically, InternVL3.5-
631 8B failed to consistently follow the required out-
632 put format unless a special `<assistant>` token was
633 removed. We did not observe this issue with
634 InternVL3.5-38b. Similarly, GPT-5.1 tended to by-
635 pass the requested chain-of-thought reasoning and
636 produce only a final answer. To fix it, we added the
637 instruction: "you must follow the exact reasoning
638 steps provided in the examples before providing the
639 final answer" at the end of the prompt. Addition-
640 ally, due to resource constraints, we were unable
641 to explore multimodal in-context learning (ICL),
642 leaving open the question of how combining image
643 and text examples in ICL would affect performance.
644 Finally, our experiments were conducted only in
645 English; thus, the findings may not generalize to
646 other languages.

647 Ethical Statements

648 Although we believe this work raises no direct eth-
649 ical concerns, we cannot fully rule out potential
650 dual-use implications of our findings.

651 References

652 Deepak Bhaskar Acharya, Karthigeyan Kuppan, and
653 B Divya. 2025. Agentic ai: Autonomous intelligence
654 for complex goals—a comprehensive survey. *IEEE*
655 *Access*.

656 James F Allen. 1983. Maintaining knowledge about
657 temporal intervals. *Communications of the ACM*,
658 26(11):832–843.

659 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-
660 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shi-
661 jie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu,
662 Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei
663 Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 oth-
664 ers. 2025. *Qwen2.5-vl technical report*. *Preprint*,
665 arXiv:2502.13923.

666 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
667 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
668 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
669 Askell, and 1 others. 2020. Language models are
670 few-shot learners. *Advances in neural information*
671 *processing systems*, 33:1877–1901.

Tom Bylander. 1994. The computational complexity of
672 propositional strips planning. *Artificial Intelligence*,
673 69(1-2):165–204. 674

Nathanael Chambers. 2013. Navytime: Event and time
675 ordering from raw text. 676

Delong Chen, Theo Moutakanni, Willy Chung, Yejin
677 Bang, Ziwei Ji, Allen Bolourchi, and Pascale Fung.
678 2025. Planning with reasoning using vision language
679 world model. *arXiv preprint arXiv:2509.02722*. 680

Leon RA Derczynski. 2017. *Automatically ordering*
681 *events and times in text*, volume 677. Springer. 682

Sorouralsadat Fatemi and Yuheng Hu. 2024. Enhancing
683 financial question answering with a multi-agent re-
684 flection framework. In *Proceedings of the 5th ACM*
685 *International Conference on AI in Finance*, pages
686 530–537. 687

Richard E Fikes and Nils J Nilsson. 1971. Strips: A new
688 approach to the application of theorem proving to
689 problem solving. *Artificial intelligence*, 2(3-4):189–
690 208. 691

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
692 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
693 Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-
694 ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh
695 Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-
696 tra, Archie Sravankumar, Artem Korenev, Arthur
697 Hinsvark, and 542 others. 2024. *The llama 3 herd of*
698 *models*. *Preprint*, arXiv:2407.21783. 699

Naresh Gupta, Dana S Nau, and 1 others. 1991. Com-
700 plexity results for blocks-world planning. In *AAAI*,
701 volume 91, pages 629–633. 702

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
703 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
704 Weizhu Chen. 2021. *Lora: Low-rank adaptation of*
705 *large language models*. *Preprint*, arXiv:2106.09685. 706

Yiyang Ji, Haoran Chen, Qiguang Chen, Chengyue Wu,
707 Libo Qin, and Wanxiang Che. 2025. *Mpcc: A novel*
708 *benchmark for multimodal planning with complex*
709 *constraints in multimodal large language models*. In
710 *Proceedings of the 33rd ACM International Confer-*
711 *ence on Multimedia*, MM '25, page 5188–5197, New
712 York, NY, USA. Association for Computing Machin-
713 ery. 714

Subbarao Kambhampati, Karthik Valmeekam, Lin
715 Guan, Mudit Verma, Kaya Stechly, Siddhant Bham-
716 bri, Lucas Saldyt, and Anil Murthy. 2024. Position:
717 LLMs can't plan, but can help planning in llm-modulo
718 frameworks. In *Proceedings of the 41st Interna-*
719 *tional Conference on Machine Learning*, ICML'24.
720 JMLR.org. 721

Yash Kumar Lal, Vanya Cohen, Nathanael Chambers,
722 Niranjana Balasubramanian, and Ray Mooney. 2024.
723 *CaT-bench: Benchmarking language model under-*
724 *standing of causal and temporal dependencies in*
725 *plans*. In *Proceedings of the 2024 Conference on*
726

727	<i>Empirical Methods in Natural Language Processing</i> , pages 19336–19354, Miami, Florida, USA. Association for Computational Linguistics.	<i>Conference on Language Resources and Evaluation (LREC'06)</i> , Genoa, Italy. European Language Resources Association (ELRA).	784
728			785
729			786
730	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. Llava-onevision: Easy visual task transfer . <i>Preprint</i> , arXiv:2408.03326.	Lu Qiu, Yi Chen, Yuying Ge, Yixiao Ge, Ying Shan, and Xihui Liu. 2025. Egoplan-bench2: A benchmark for multimodal large language model planning in real-world scenarios . <i>Preprint</i> , arXiv:2412.04447.	787
731			788
732			789
733			790
734			
735	Kanxue Li, Baosheng Yu, Qi Zheng, Yibing Zhan, Yuhui Zhang, Tianle Zhang, Yijun Yang, Yue Chen, Lei Sun, Qiong Cao, Li Shen, Lusong Li, Dapeng Tao, and Xiaodong He. 2024b. Muep: A multimodal benchmark for embodied planning with foundation models . In <i>Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24</i> , pages 129–138. International Joint Conferences on Artificial Intelligence Organization. Main Track.	Tayiba Raheem and Gahangir Hossain. 2025. Agentic ai systems: Opportunities, challenges, and trustworthiness. In <i>2025 IEEE International Conference on Electro Information Technology (eIT)</i> , pages 618–624. IEEE.	791
736			792
737			793
738			794
739			795
740			
741		Anand S Rao, Michael P Georgeff, and 1 others. 1995. Bdi agents: From theory to practice. In <i>Icmas</i> , volume 95, pages 312–319.	796
742			797
743			798
744			
745	Fangru Lin, Emanuele La Malfa, Valentin Hofmann, Elle Michelle Yang, Anthony G. Cohn, and Janet B. Pierrehumbert. 2024. Graph-enhanced large language models in asynchronous plan reasoning. In <i>Proceedings of the 41st International Conference on Machine Learning, ICML'24</i> . JMLR.org.	Gabriel Roccabruna, Massimo Rizzoli, Giuseppe Riccardi, and 1 others. 2024. Will llms replace the encoder-only models in temporal relation classification? In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	799
746			800
747			801
748			802
749			803
750			804
751	Yujie Lu, Pan Lu, Zhiyu Chen, Wanrong Zhu, Xin Eric Wang, and William Yang Wang. 2024. Multimodal procedural planning via dual text-image prompting . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10931–10954, Miami, Florida, USA. Association for Computational Linguistics.	Samdyuti Suri, Sankar Narayan Das, Kapil Singi, Kuntal Dey, Vibhu Saujanya Sharma, and Vikrant Kaulgud. 2023. Software engineering using autonomous agents: Are we there yet? In <i>2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)</i> , pages 1855–1857. IEEE.	805
752			806
753			807
754			808
755			809
756			810
757			811
758	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, and 1 others. 2023. Self-refine: Iterative refinement with self-feedback. <i>Advances in Neural Information Processing Systems</i> , 36:46534–46594.	Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. 2023. On the planning abilities of large language models: a critical investigation. In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23</i> , Red Hook, NY, USA. Curran Associates Inc.	812
759			813
760			814
761			815
762			816
763			817
764	Tianjun Mo, Ziyi Jiang, and Qichang Zheng. 2025. Interactive ai agent for code refactoring assistance: A study on decision-making strategies and human-agent collaboration effectiveness. <i>Academia Nexus Journal</i> , 4(1).	Siddharth Vashishtha, Adam Poliak, Yash Kumar Lal, Benjamin Van Durme, and Aaron Steven White. 2020. Temporal reasoning in natural language inference. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4070–4078.	819
765			820
766			821
767			822
768			823
769	Shinsuke Mori, Hirokuni Maeta, Yoko Yamakata, and Tetsuro Sasada. 2014. Flow graph corpus from recipe texts . In <i>Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)</i> , pages 2370–2377, Reykjavik, Iceland. European Language Resources Association (ELRA).	Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, Zhaokai Wang, Zhe Chen, Hongjie Zhang, Ganlin Yang, Haomin Wang, Qi Wei, Jinhui Yin, Wenhao Li, Erfei Cui, and 56 others. 2025. Internvl3.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency . <i>Preprint</i> , arXiv:2508.18265.	824
770			825
771			826
772			827
773			828
774			829
775	Liangming Pan, Jingjing Chen, Jianlong Wu, Shaoteng Liu, Chong-Wah Ngo, Min-Yen Kan, Yu-Gang Jiang, and Tat seng Chua. 2020. Multi-modal cooking workflow construction for food recipes . <i>Proceedings of the 28th ACM International Conference on Multimedia</i> .	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	830
776			831
777			832
778			833
779			834
780			835
781	Rebecca Passonneau. 2006. Measuring agreement on set-valued items (MASI) for semantic and pragmatic annotation . In <i>Proceedings of the Fifth International</i>		836
782			837
783			

838 Joseph Weizenbaum. 1966. Eliza—a computer program
839 for the study of natural language communication be-
840 tween man and machine. *Communications of the*
841 *ACM*, 9(1):36–45.

842 Michael Wooldridge and Nicholas R Jennings. 1995. In-
843 telligent agents: Theory and practice. *The knowledge*
844 *engineering review*, 10(2):115–152.

845 Te-Lin Wu, Alex Spangher, Pegah Alipoormolabashi,
846 Marjorie Freedman, Ralph Weischedel, and Nanyun
847 Peng. 2022. [Understanding multimodal procedural](#)
848 [knowledge by sequencing multimodal instructional](#)
849 [manuals](#). In *Proceedings of the 60th Annual Meet-*
850 *ing of the Association for Computational Linguistics*
851 *(Volume 1: Long Papers)*, pages 4525–4542, Dublin,
852 Ireland. Association for Computational Linguistics.

853 Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze
854 Lou, Yuandong Tian, Yanghua Xiao, and Yu Su.
855 2024. Travelplanner: a benchmark for real-world
856 planning with language agents. In *Proceedings of the*
857 *41st International Conference on Machine Learning*,
858 ICML’24. JMLR.org.

859 Yoko Yamakata, Shinsuke Mori, and John Carroll. 2020.
860 [English recipe flow graph corpus](#). In *Proceedings*
861 *of the Twelfth Language Resources and Evaluation*
862 *Conference*, pages 5187–5194, Marseille, France. Eu-
863 ropean Language Resources Association.

864 Wentao Zhang, Lingxuan Zhao, Haochong Xia, Shuo
865 Sun, Jiaze Sun, Molei Qin, Xinyi Li, Yuqing Zhao,
866 Yilei Zhao, Xinyu Cai, and 1 others. 2024. A multi-
867 modal foundation agent for financial trading: Tool-
868 augmented, diversified, and generalist. In *Proceed-*
869 *ings of the 30th acm sigkdd conference on knowledge*
870 *discovery and data mining*, pages 4314–4325.

871 Xi Zhang and Qixuan Zhu. 2023. Ai-enabled network-
872 functions virtualization and software-defined archi-
873 tectures for customized statistical qos over 6g mas-
874 sive mimo mobile wireless networks. *IEEE Network*,
875 37(2):30–37.

876 Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang,
877 Xinyun Chen, Minmin Chen, Azade Nova, Le Hou,
878 Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, and
879 Denny Zhou. 2024. [Natural plan: Benchmark-](#)
880 [ing llms on natural language planning](#). *Preprint*,
881 arXiv:2406.04520.

890 7b¹³, InternVL3.5-8b¹⁴, InternVL3.5-38b¹⁵, and
891 LLaVA-OneVision-7b¹⁶.
892 To run closed-source GPT-5.1¹⁷ model, we used
893 OpenAI Batch API¹⁸. As with the open-source
894 models, we used greedy decoding for reproducibil-
895 ity (temperature=0). Since we wanted the model
896 to follow our reasoning steps (e.g. CoT exper-
897 iment), we disabled the model’s default inter-
898 nal reasoning by setting reasoning={"effort":
899 "none"}.

900 For fine-tuning, we used a batch size of 1, a
901 learning rate of 1e-5, and a weight decay of 0.01.
902 The models were fine-tuned for 1 epoch. We used
903 Low-Rank Adaptation (LoRA) (Hu et al., 2021) to
904 fine-tune the model. We set alpha to 32 and rank
905 to 16. LoRA adapters are applied to the query, key,
906 value, and output matrices of the attention layers.

907 Experiments involving Qwen2.5-VL-72B,
908 InternVL3.5-38B, and fine-tuning were conducted
909 using two NVIDIA A100 GPUs with 80 GiB of
910 memory, while all other experiments were run on
911 a single NVIDIA A100 with 80 GiB of memory.
912 For the experiments with GPT-5.1, we spent a total
913 of around \$100, considering both experimentation
914 and computing the results presented in this article.
915 While for collecting the corpus, we spent a total of
916 around £800, considering both tests and the actual
917 annotation.

918 The parameters used to screen the participant
919 sample provided by Prolific were as follows:

- First language: Italian 920
- Country of Birth: Italy 921
- Education: High school diploma/A-levels or
above 922
- Approval rate: 95-100 924
- Number of previous submissions: 100-10000 925

A.2 Additional Material 926

A Appendix

A.1 Experimental details

We used the HuggingFace framework to run experiments with open-source models. For generation, we employed greedy decoding (do_sample=False, num_beams=1) for all open-source models to foster reproducibility. We used the following model checkpoints: Qwen2.5-VL-72b¹², Qwen2.5-VL-

¹²<https://huggingface.co/Qwen/Qwen2.5-VL-72B-Instruct>

¹³<https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct>

¹⁴https://huggingface.co/OpenGVLab/InternVL3_5-8B-HF

¹⁵https://huggingface.co/OpenGVLab/InternVL3_5-38B-HF

¹⁶<https://huggingface.co/lmms-lab/llava-onevision-qwen2-7b-ov>

¹⁷Model snapshot: gpt-5.1-2025-11-13

¹⁸<https://platform.openai.com/docs/guides/batch>

Prompt Setting	Modality	Prompt
Baseline	Image-only	Using ONLY the information in the Context, answer the following three questions in EXACTLY this format: Q1: The answer is: <Yes/No/I don't know>. Q2: The answer is: <Yes/No/I don't know>. Q3: The answer is: <Yes/No/I don't know>. Do not add anything else. Do not explain. Do not change the format. Context: Step A picture: [image_A] Step B picture: [image_B] Questions: Q1: Must Step A be executed before Step B? Q2: Must Step A be executed after Step B? Q3: Can Step A and Step B be executed in parallel?
	Text-only	Using ONLY the information in the Context ... Context: Step A description: [description_A] Step B description: [description_B] Questions ...
	Image-Text	Using ONLY the information in the Context ... Context: Step A picture: [image_A] Step A description: [description_A] Step B picture: [image_B] Step B description: [description_B] Questions ...
Instructions	Image-only	Your task is to determine the dependency order between two steps in a recipe. Follow these rules: - Before : Step A must be executed before Step B if the outcome of Step A is required to complete Step B (i.e., Step B depends on Step A). - After : Step A must be executed after Step B if the outcome of Step B is required to complete Step A (i.e., Step A depends on Step B). - Parallel : Step A and Step B can be executed in parallel if neither step depends on the outcome of the other; therefore, their order of execution can be arbitrary. Using ONLY the information in the Context ... Context: ... Questions: ...
	Text-only	Your task is to determine the dependency order between two steps in a recipe. Follow these rules: - Before: ... - After: ... - Parallel: ... Ignore sequencing terms (e.g., 'first', 'then', 'lastly', and other words that may appear in the text for the natural flow of the recipe) when determining the execution order, and focus only on the action itself.

Prompt Setting	Modality	Prompt
		Using ONLY the information in the Context ... Context: ... Questions: ...
	Image-Text	Your task is to determine the dependency order between two steps in a recipe. Follow these rules: - Before: ... - After: ... - Parallel: ... Ignore sequencing terms (e.g., 'first', 'then', 'lastly', and other words that may appear in the text for the natural flow of the recipe) when determining the execution order, and focus only on the action itself. Also note that the text description may include partial or full references to steps not shown in the image; in such cases, rely on the actions depicted in the image. Using ONLY the information in the Context ... Context: ... Questions: ...
In-context learning	Image-only	Your task is to determine the dependency order between two steps in a recipe. Follow these rules: - Before: ... - After: ... - Parallel: ... You will be shown three examples demonstrating how to solve the task using text-based step descriptions. However, your actual input will consist of images, and your reasoning should be based on the actions depicted in those images. Examples: Step A description: Grate the lemon zest. Step B description: Put the grated lemon zest into the strawberry sauce. Questions: Q1: Must Step A be executed before Step B? Q2: Must Step A be executed after Step B? Q3: Can Step A and Step B be executed in parallel? Q1: The answer is: Yes. Q2: The answer is: No. Q3: The answer is: No. Explanation: Step B explicitly depends on Step A - lemon zest must already be grated (Step A) before it can be put into the strawberry sauce (Step B); therefore, Step A must be executed before Step B. Step A description: Add the celery. Step B description: Then add carrots. Questions: Q1: Must Step A be executed before Step B? Q2: Must Step A be executed after Step B? Q3: Can Step A and Step B be executed in parallel? Q1: The answer is: No.

Prompt Setting	Modality	Prompt
		<p>Q2: The answer is: No. Q3: The answer is: Yes. Explanation: Both actions are independent; neither step produces something the other one requires.</p> <p>Step A description: Pour it into the cup. Step B description: Measure out raspberry juice. Questions: Q1: Must Step A be executed before Step B? Q2: Must Step A be executed after Step B? Q3: Can Step A and Step B be executed in parallel? Q1: The answer is: No. Q2: The answer is: Yes. Q3: The answer is: No. Explanation: Step A relies on the outcome of Step B - raspberry juice must be poured into the cup (Step A) after it is measured out (Step B); therefore, Step A must be executed after Step B. Using ONLY the information in the Context ... Context: ... Questions: ...</p>
	Text-only	<p>Your task is to determine the dependency order between two steps in a recipe. Follow these rules:</p> <ul style="list-style-type: none"> - Before: ... - After: ... - Parallel: ... <p>Ignore sequencing terms ... You will be shown three examples demonstrating how to solve the task using text-based step descriptions. Examples: ... Using ONLY the information in the Context ... Context: ... Questions: ...</p>
	Image-Text	<p>Your task is to determine the dependency order between two steps in a recipe. Follow these rules:</p> <ul style="list-style-type: none"> - Before: ... - After: ... - Parallel: ... <p>Ignore sequencing terms ... You will be shown three examples demonstrating how to solve the task using text-based step descriptions. However, your actual input will consist of both images and text descriptions, and your reasoning should be based on both actions shown in the images and the accompanying textual descriptions. Examples: ... Using ONLY the information in the Context ... Context: ... Questions: ...</p>
CoT	Image-only	<p>Your task is to determine the dependency order between two steps in a recipe. You must choose from: Before, After, or Parallel. Follow these rules:</p>

Prompt Setting	Modality	Prompt
		<p>- Before: ... - After: ... - Parallel: ...</p> <p>You will be shown three examples demonstrating how to solve the task using text-based step descriptions. However, your actual input will consist of images, and your reasoning should be based on the actions depicted in those images. You must follow the reasoning steps shown in the examples before answering.</p> <p>Examples:</p> <p>Step A description: Grate the lemon zest. Step B description: Put the grated lemon zest into the strawberry sauce. Step A produces: Grated lemon zest. Step B produces: Lemon zest inside the strawberry sauce. Step A requires: A lemon. Step B requires: Lemon zest that has been grated. Dependency analysis: Step B explicitly depends on Step A - lemon zest must already be grated (Step A) before it can be put into the strawberry sauce (Step B); therefore, Step A must be executed before Step B. The answer is: Before.</p> <p>Step A description: Add the celery. Step B description: Then add carrots. Step A produces: A component with the celery added. Step B produces: A component with the carrots added. Step A requires: The celery. Step B requires: The carrots. Dependency analysis: Each step adds a separate ingredient, and neither depends on the other, so they can occur in any order. The answer is: Parallel.</p> <p>Step A description: Pour it into the cup. Step B description: Measure out raspberry juice. Step A produces: The cup with raspberry juice poured in it. Step B produces: Raspberry juice that was measured out. Step A requires: Raspberry juice that was measured out. Step B requires: Raspberry juice. Dependency analysis: Step A relies on the outcome of Step B - raspberry juice must be poured into the cup (Step A) after it is measured out (Step B); therefore, Step A must be executed after Step B. The answer is: After.</p> <p>Step A picture: ...</p>
Text-only		<p>Your task is to determine the dependency order between two steps in a recipe. You must choose from: Before, After, or Parallel. Follow these rules:</p>

Prompt Setting	Modality	Prompt
		<p>- Before: ... - After: ... - Parallel: ... Ignore sequencing terms ... You will be shown three examples demonstrating how to solve the task using text-based step descriptions. You must follow the reasoning steps shown in the examples before answering.</p> <p>Examples: ...</p>
		<p>Step A description: ...</p>
	Image-Text	<p>Your task is to determine the dependency order between two steps in a recipe. You must choose from: Before, After, or Parallel. Follow these rules:</p> <p>- Before: ... - After: ... - Parallel: ...</p> <p>You will be shown three examples demonstrating how to solve the task using text-based step descriptions. However, your actual input will consist of both images and text descriptions, and your reasoning should be based on both actions shown in the images and the accompanying textual descriptions. You must follow the reasoning steps shown in the examples before answering.</p> <p>Examples: ...</p> <p>Step A picture: ... Step A description: ...</p>
Self-reflection	Image-only	<p>Your task is to determine the dependency order between two steps in a recipe. You must choose from: Before, After, or Parallel. Follow these rules:</p> <p>- Before: ... - After: ... - Parallel: ...</p> <p>You will be shown three examples demonstrating how to solve the task using text-based step descriptions. However, your actual input will consist of images, and your reasoning should be based on the actions depicted in those images. You must follow the reasoning steps shown in the examples before answering. After you answer the question, review your reasoning and check whether your answer logically follows from the context and dependencies you identified. After self-reflection, provide your final answer, confirming or correcting your initial choice.</p> <p>Examples:</p> <p>Step A description: Grate the lemon zest. Step B description: Put the grated lemon zest into the strawberry sauce.</p>

Prompt Setting	Modality	Prompt
		<p>Step A produces: Grated lemon zest. Step B produces: Lemon zest inside the strawberry sauce. Step A requires: A lemon. Step B requires: Lemon zest that has been grated. Dependency analysis: Step B explicitly depends on Step A - lemon zest must already be grated (Step A) before it can be put into the strawberry sauce (Step B); therefore, Step A must be executed before Step B. The answer is: Before. Reflection: <your_reflection> The final answer: <final_answer></p> <p>...</p>
	Text-only	<p>Step A picture: ...</p> <p>Your task is to determine the dependency order between two steps in a recipe. You must choose from: Before, After, or Parallel. Follow these rules:</p> <ul style="list-style-type: none"> - Before: ... - After: ... - Parallel: ... <p>Ignore sequencing terms ...</p> <p>You will be shown three examples demonstrating how to solve the task using text-based step descriptions. You must follow the reasoning steps shown in the examples before answering. After you answer the question, review your reasoning and check whether your answer logically follows from the context and dependencies you identified. After self-reflection, provide your final answer, confirming or correcting your initial choice.</p> <p>Examples: ...</p>
	Image-Text	<p>Step A description: ...</p> <p>Your task is to determine the dependency order between two steps in a recipe. You must choose from: Before, After, or Parallel. Follow these rules:</p> <ul style="list-style-type: none"> - Before: ... - After: ... - Parallel: ... <p>You will be shown three examples demonstrating how to solve the task using text-based step descriptions. However, your actual input will consist of both images and text descriptions, and your reasoning should be based on both actions shown in the images and the accompanying textual descriptions. You must follow the reasoning steps shown in the examples before answering.</p> <p>Examples: ...</p> <p>Step A picture: ...</p>

Prompt Setting	Modality	Prompt
		Step A description: ...

Table 5: Prompt schema used for models evaluation under image-only, text-only, and image-text modalities. Repeated or previously described components are omitted for brevity.

Modality	Prompt
Image-only	<p data-bbox="424 275 1342 338"><User> Your task is to determine the dependency order between two steps in a recipe. Follow these rules:</p> <ul data-bbox="424 344 1342 629" style="list-style-type: none"> <li data-bbox="424 344 1342 407">- Before: Step A must be executed before Step B if the outcome of Step A is required to complete Step B (i.e., Step B depends on Step A). <li data-bbox="424 414 1342 477">- After: Step A must be executed after Step B if the outcome of Step B is required to complete Step A (i.e., Step A depends on Step B). <li data-bbox="424 483 1342 546">- Parallel: Step A and Step B can be executed in parallel if neither step depends on the outcome of the other; therefore, their order of execution can be arbitrary. <p data-bbox="424 553 735 582">Step A picture: [image_A]</p> <p data-bbox="424 589 735 618">Step B picture: [image_B]</p> <p data-bbox="424 624 549 654">Questions:</p> <p data-bbox="424 660 951 689">Q1: Must Step A be executed before Step B?</p> <p data-bbox="424 696 930 725">Q2: Must Step A be executed after Step B?</p> <p data-bbox="424 732 1031 761">Q3: Can Step A and Step B be executed in parallel?</p> <p data-bbox="424 768 564 797"><Assistant></p> <p data-bbox="424 804 703 833">Q1: The answer is: Yes.</p> <p data-bbox="424 840 699 869">Q2: The answer is: No.</p> <p data-bbox="424 875 699 904">Q3: The answer is: No.</p>
Text-only	<p data-bbox="424 936 1342 999"><User> Your task is to determine the dependency order between two steps in a recipe. Follow these rules:</p> <ul data-bbox="424 1005 1342 1144" style="list-style-type: none"> <li data-bbox="424 1005 560 1034">- Before: ... <p data-bbox="424 1041 1342 1144">Ignore sequencing terms (e.g., 'first', 'then', 'lastly', and other words that may appear in the text for the natural flow of the recipe) when determining the execution order, and focus only on the action itself.</p> <p data-bbox="424 1151 842 1180">Step A description: [description_A]</p> <p data-bbox="424 1187 842 1216">Step B description: [description_B]</p> <p data-bbox="424 1223 549 1252">Questions:</p> <p data-bbox="424 1258 951 1288">Q1: Must Step A be executed before Step B?</p> <p data-bbox="424 1294 930 1323">Q2: Must Step A be executed after Step B?</p> <p data-bbox="424 1330 1031 1359">Q3: Can Step A and Step B be executed in parallel?</p> <p data-bbox="424 1366 564 1395"><Assistant></p> <p data-bbox="424 1402 699 1431">Q1: The answer is: No.</p> <p data-bbox="424 1438 699 1467">Q2: The answer is: No.</p> <p data-bbox="424 1473 703 1503">Q3: The answer is: Yes.</p>
Image-Text	<p data-bbox="424 1525 1342 1588"><User> Your task is to determine the dependency order between two steps in a recipe. Follow these rules:</p> <ul data-bbox="424 1594 1342 1809" style="list-style-type: none"> <li data-bbox="424 1594 560 1624">- Before: ... <p data-bbox="424 1637 1342 1809">Ignore sequencing terms (e.g., 'first', 'then', 'lastly', and other words that may appear in the text for the natural flow of the recipe) when determining the execution order, and focus only on the action itself. Also note that the text description may include partial or full references to steps not shown in the image; in such cases, rely on the actions depicted in the image.</p> <p data-bbox="424 1816 842 1845">Step A description: [description_A]</p> <p data-bbox="424 1852 842 1881">Step B description: [description_B]</p> <p data-bbox="424 1888 549 1917">Questions:</p> <p data-bbox="424 1924 951 1953">Q1: Must Step A be executed before Step B?</p> <p data-bbox="424 1960 930 1989">Q2: Must Step A be executed after Step B?</p> <p data-bbox="424 1995 1031 2024">Q3: Can Step A and Step B be executed in parallel?</p> <p data-bbox="424 2031 564 2060"><Assistant></p>

Q1: The answer is: No.
Q2: The answer is: Yes.
Q3: The answer is: No.

Table 6: Example of prompt–response pairs used for fine-tuning across image-only, text-only, and image–text modalities. Each example includes instructions, input structure, and expected binary outputs. The <User> and <Assistant> tokens are model-specific and are included here solely for illustrative purposes.

1. Purpose

Welcome! The goal of our research project is to test whether humans can recognize the order of steps in a recipe. For this, we need annotated data, so we need your help!

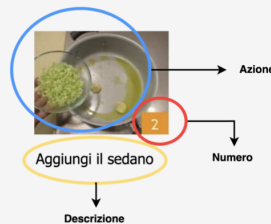
Note

The assignment will be divided into two parts. The first part includes a test (**Comprehension Check**) to verify your understanding of the assignment. If you pass, you'll move on to the second part, where you'll have to write down a series of recipes.

In the test part you will have to write down 2 recipes and in the annotation part you will have to write down 5 recipes.

2. Task

In this study, we'll ask you to note **the order of dependencies between the steps of a recipe** . We'll provide you with all the steps of a recipe along with the corresponding images showing: the action , the text description, and the step number that reflects the order of the original recipe. Note that the step number is only used to follow the flow of the recipe and may **not** correlate with the order of dependencies.



For each step of a recipe:

- **Find the part of the description related to the image:** first, you'll need to read the description and find the part of the text that best describes or relates to the action depicted in the image. *Use this part of the description to decide the order of dependencies!* The rest of the text will help you understand the flow of the recipe and should therefore be used as context.
- **Otherwise, rely on the image:** in some cases, the description of the action may not be present in the text. In these cases, we ask you to rely on the action depicted in the image, or what the image communicates.

For clarity, in the examples below, we've removed the passage description from the image, but you'll see it displayed in the annotation activity.

Figure 2: This figure shows the first two points of the annotation guidelines, that present the purpose and the general description of the task.

A.3 AI tool usage

During the development phase, we used ChatGPT and Gemini to identify and resolve implementation bugs. We also leveraged these models and Grammarly to suggest rephrasing that improves the fluency and readability of the manuscript. Additionally, we have experimented with using Google Scholar Labs to expand and refine our literature review.

A.4 Licences of the corpus and models used

All models used in our experiments have been utilised in accordance with their respective licenses. We plan to license the corpus with the Creative Commons 4.0 CC BY-NC-SA 4.0 (Attribution-NonCommercial-ShareAlike). The intended use of this corpus will be only for scientific research purposes.

3. Order of dependency

To identify the dependency order **watch the video** below and **read the descriptions** 1) dependency order, 2) concurrent execution and 3) transitivity in **green** and the related examples in **blue**.

How to identify the dependency order?



1) Order of Dependence

Given two steps **J** and **P**, **P** is required by **J** ($P \rightarrow J$) if step **P** needs to be completed before step **J**. In other words, **the result of step P is absolutely necessary for step J to be completed**. You can think of steps as processes. A process needs some requirements to begin with. If those requirements are produced by another process, then there is a dependency between the two. To help understand better, you can refer to Example 1.1.

Example 1.1

Given the steps:

11 Grate the lemon zest;

12 Add the grated lemon zest to the strawberry sauce;

12 depends on **11**. This is because without grating the lemon zest, it is not possible to add the grated zest to the strawberries.



Figure 3: This section of the annotation guidelines offers annotators the option to watch a demonstration video that walks through the task and explains how to identify dependencies between steps. Following the video, annotators are provided with a textual description that formally defines each dependency, accompanied by illustrative examples.

Example 1.1

Given the steps:

11 Grate the lemon zest;

12 Add the grated lemon zest to the strawberry sauce;

12 depends on **11**. This is because without grating the lemon zest, it is not possible to add the grated zest to the strawberries.



2) Simultaneous execution

There may be steps that have no dependencies between them, meaning they can theoretically be executed simultaneously. That is, the order in which they are executed does not affect the correct execution of the recipe. To help you understand, follow Examples 2.1, 2.2, and 2.3.

Example 2.1

Given the steps:

1 Fry a clove of garlic with a drizzle of extra virgin olive oil;

2 Add the celery;

3 Add carrots;

4 Add onions;

Step **1** must be done first. This is because the garlic must be sautéed before proceeding to the next steps. Steps **2**, **3**, and **4** can be done simultaneously. That is **1** → **2**, **1** → **3**, and **1** → **4**, as shown in the figure. This is because I don't need onions to correctly perform the "add carrots" action, or vice versa. Furthermore, the recipe doesn't explicitly specify any order, so we can be sure that there is no order dependency between these steps.

Figure 4: This figure continues from the previous one, illustrating a case in which a step depends on another (Example 1.1), and presenting the formal definition of independent steps

2) Simultaneous execution

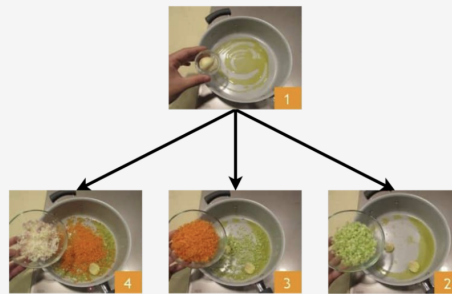
There may be steps that have no dependencies between them, meaning they can theoretically be executed simultaneously. That is, the order in which they are executed does not affect the correct execution of the recipe. To help you understand, follow Examples 2.1, 2.2, and 2.3.

Example 2.1

Given the steps:

- 1 Fry a clove of garlic with a drizzle of extra virgin olive oil;
- 2 Add the celery;
- 3 Add carrots;
- 4 Add onions;

Step 1 must be done first. This is because the garlic must be sautéed before proceeding to the next steps. Steps 2, 3, and 4 can be done simultaneously. That is 1 → 2, 1 → 3, and 1 → 4, as shown in the figure. This is because I don't need onions to correctly perform the "add carrots" action, or vice versa. Furthermore, the recipe doesn't explicitly specify any order, so we can be sure that there is no order dependency between these steps.



Example 2.2

The part that follows Example 2.1 is that steps that can ideally be performed simultaneously must finish before step 5. So it follows that 5 depends on 2, 3, and 4. The reason is that you can't start with step 5 without having the celery, carrots, and onions in the pan.



Figure 5: This figure shows an example given to illustrate the case in which steps are independent.

Example 2.3

Given the following steps:

- 1 Add ice to a glass;
- 2 Measure out the lemon juice;
- 3 Pour it;
- 4 Measure out the grape juice;
- 5 Pour it;
- 6 Measure and pour in the raspberry juice;
- 7 Mix well;

We can note that steps 2 and 3, and 4 and 5 are the breakdown of step 6. Furthermore, considering that in step 7 everything is mixed, that the recipe does not specify the order of additions, and that, for example, I do not need the lemon juice to physically pour the grape juice, we can say with good certainty that the pairs 2 - 3, 4 - 5 and step 6 can be performed simultaneously.

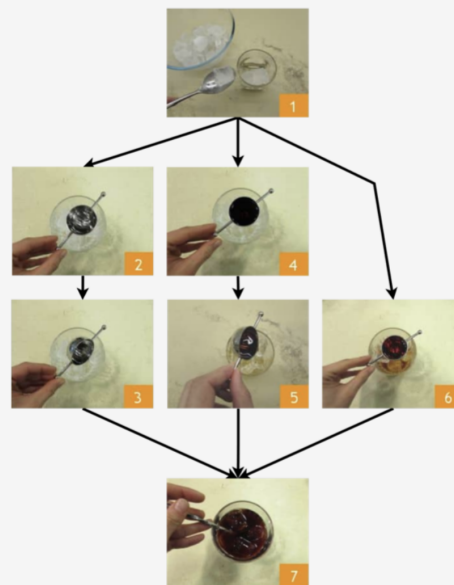
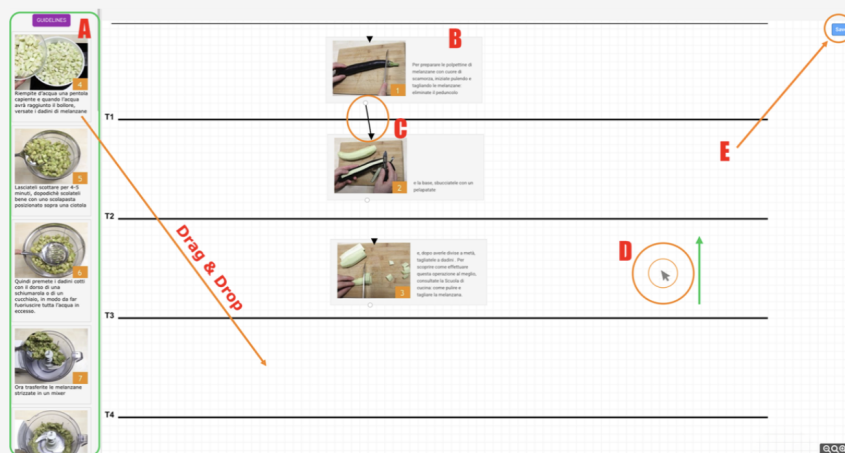


Figure 6: Example illustrating a fully annotated DAG for a recipe, with all relations clearly explained.

5. Annotation tool

You should now have a clear understanding of the task. So, let's briefly review the annotation tool. The annotation tool is illustrated in the following image:



- **A:** This is the panel with the list of recipe steps you can add to the timeline. At the top of this panel, you'll find a "Guidelines" button. This button takes you back to the guidelines you're currently reading, so you can refer back to them at any time during your annotation (for example, if you're unsure or forget something).
- **B:** To drag steps onto the timeline, you need to hold down the left mouse button on the gray area of the step.
- **C:** To connect two steps, click the white dot of the first step and then the black triangle of the second step. This will create a line showing the dependency between the two steps. To delete a connection, simply click the connection line you want to delete and press "delete" or, on the keyboard, press *DEL*, *DEL*, or on a MAC, *Cmd+Delete/Backspace*.
- **I:** To scroll down the timeline, hold down the left mouse button on any part of the timeline and move the mouse up. (*It's like scrolling a web page on a smartphone.*) The same goes for moving left and right.
- **E:** once you have connected all the steps, click on the "Save" button to proceed to the next recipe or finish the activity, in case there are no more recipes.
 - **Note:** There may be some rare cases where a step can't be linked to anything. In this case, once you've finished linking all the other steps, click "Save," and we'll then ask you to provide the reason why this step couldn't be linked to the others.

Figure 7: Illustration of the annotation tool as presented in the annotation guidelines. Highlighted regions (marked in red letters) indicate key functionalities, including dragging steps, adding or deleting connections, and saving the annotations.



Figure 8: Example of a DAG annotated by a crowd-worker. Steps 1, 2, 3, 4, and 5 can be executed in any order (or simultaneously), thus there is no directed edge between them.