

Conversational Goal-Conflict Explanations in Planning via Multi-Agent LLMs

Submission #4858

Abstract

1 When automating plan generation for a real-world
2 sequential decision problem, the goal is often not
3 to replace the human planner, but to facilitate an
4 iterative reasoning and elicitation process, where
5 the human’s role is to guide the planner accord-
6 ing to their preferences and expertise. In this con-
7 text, explanations that respond to users’ ques-
8 tions are crucial to improve their understanding of po-
9 tential solutions and increase their trust in the system.
10 To enable natural interaction with such a system,
11 we present a multi-agent Large Language Model
12 (LLM) architecture that is agnostic to the expla-
13 nation framework and enables user- and context-
14 dependent interactive explanations. We also de-
15 scribe an instantiation of this framework for goal-
16 conflict explanations, which we use to conduct a
17 pilot case study comparing the LLM-powered in-
18 teraction with a more conventional template-based
19 explanation interface.

20 1 Introduction

21 One can think of planning as the task of finding a plan that sat-
22 isfies a set of properties, but also as the iterative process that
23 starts before the goals, objectives and preferences are fully
24 defined [Smith, 2012], and ends when a plan is found that
25 is satisfactory for all parties involved. From this perspective, ex-
26 planations serve the purpose of accelerating the convergence
27 of preferences elicitation by humans. Interactive planning has
28 become a relatively well-accepted paradigm, yet the crucial
29 role of interactive *explanations* in that process has been less
30 widely investigated. Explanations are by nature interactive,
31 as explanation is something that one person (the explainer)
32 does for the sake of another (the explainee). The request
33 for an explanation arises because the explainee has a concep-
34 tual problem [Bromberger, 1962; Achinstein, 1980]; there is
35 something they do not understand. The explainer’s task is to
36 resolve the conceptual problem. However, fixing a concep-
37 tual problem may not be a one shot deal; it may take several
38 interactions to arrive at an explanation that is satisfactory for
39 the explainee – hence the inherent interactivity of explana-
40 tions. Furthermore, previous work in psychology and human-
41 machine interaction [Liao *et al.*, 2020; Dazeley *et al.*, 2021;
42 Lakkaraju *et al.*, 2022; Zhang *et al.*, 2024] argue that in-
43 teractive explanations are more interpretable and effective

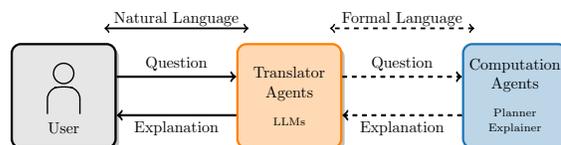


Figure 1: Architecture of our approach to interactive planning with explanations. Translator agents based on LLMs translate the user input into the formal language required by the computation agents.

44 than static ones [Miller, 2019]. Other methods were pro-
45 posed for interactive explanations [Nguyen *et al.*, 2023;
46 Shen *et al.*, 2023; Slack *et al.*, 2023; State *et al.*, 2023;
47 Feldhus *et al.*, 2023]; they use a similar high level architec-
48 ture as Figure 1, but none of them are designed for planning
49 systems.

50 If explanations are ideally interactive, then to provide ap-
51 propriate explanations, we need something with conversa-
52 tional capacities, something that is able to respond to an ini-
53 tial request for an explanation, but also to follow-up ques-
54 tions and remarks. LLMs fill this need; they are conversa-
55 tionally fluent and have impressive flexibility in translating
56 natural language to formal languages, which a symbolic plan-
57 ner even when coupled with a template interface cannot do.
58 However, LLMs cannot do everything. Empirical evidence
59 demonstrates that LLMs cannot reliably perform complex
60 planning tasks independently [Valmeekam *et al.*, 2023] and
61 there are theoretical reasons to suspect models using current
62 transformer architectures never will [Peng *et al.*, 2024]. Nev-
63 ertheless, they can serve valuable supporting roles [Kamb-
64 hampati *et al.*, 2024]. Our work extends this view as depicted
65 in Figure 1, showing how LLMs can be effectively integrated
66 into the planning process – specifically for explanations and
67 natural language interaction – while delegating the core com-
68 putational planning tasks to specialized algorithms designed
69 for that purpose.

70 To enable natural interactions between users and planning
71 systems, we need to accommodate different ways in which
72 a user formulates their question. To this end, we leverage
73 LLMs in two distinct roles: as a classifier of the question type
74 and as a translator of the question topic. Explanations should
75 resolve a conceptual problem of the explainee, which is typ-
76 ically linguistically introduced by certain types of questions.
77 These types determine the form of the explanation; answers
78 to *why* questions differ from answers to a *what-if* question;
79 e.g., “*Why can’t I go to my fitness course?*” versus “*What*”

80 *if I go grocery shopping before the fitness course?*". How- 137
 81 ever, the type of question by itself does not suffice to provide 138
 82 an explanation; we also need to specify the topic that the ex- 139
 83 planation must address. For example, a question like "*Why* 140
 84 *can't I go to my fitness course?*" has a different answer from 141
 85 "*Why can't I visit my friend?*". The explanation framework 142
 86 responds to the question-type and the question topic, trans-
 87 lated into a formal specification, with a formal explanation.
 88 To deliver this explanation to the user, we use an LLM as an
 89 explanation translator. It converts this formal explanation into
 90 a natural language response that can serve as the basis of an
 91 extended dialogue in which the original explanation is refined
 92 to meet the needs of the user.

93 The paper is organized as follows. Section 2 introduces
 94 the necessary background. Section 3 presents our general
 95 framework for iterative planning with explanations, and de-
 96 tails where translation between natural and formal language
 97 is required. Section 4 provides an implementation of expla-
 98 nations that addresses common user questions in this frame-
 99 work, whereas Section 5 describes our implementation of
 100 the LLM-based translators. Finally, Section 6 compares the
 101 LLM-based with a template-based version of the framework
 102 in a pilot user study.

103 2 Background

104 2.1 Large Language Models

105 Large Language Models (LLMs) are neural networks based
 106 on the transformer architecture [Vaswani *et al.*, 2017] that
 107 have been pre-trained on vast amounts of text data. These
 108 models process text as sequences of tokens and use self-
 109 attention mechanisms to capture relationships between dif-
 110 ferent parts of the input. Through their pre-training, LLMs
 111 have acquired capabilities in various natural language tasks,
 112 including text generation, translation, summarization and un-
 113 derstanding complex instructions. Such models are typi-
 114 cally referred to as foundation models [Touvron *et al.*, 2023;
 115 Bommasani *et al.*, 2021] when it is possible to access their
 116 parameters to specialize them to a specific purpose.

117 LLMs can be adapted in two main ways: fine-tuning,
 118 which involves additional training on specific tasks, or in-
 119 context learning, which uses carefully crafted prompts to
 120 guide the model's behavior [Brown *et al.*, 2020]. Here we use
 121 the latter approach, which has been widely adopted since it
 122 allows users to adapt LLM behavior without parameter mod-
 123 ifications, often requiring fewer than 10 examples to achieve
 124 strong performance across many tasks.

125 Multi-Agent LLM Approaches

126 Multi-agent LLM approaches involve multiple large language
 127 models working collaboratively to solve complex tasks [Guo
 128 *et al.*, 2024]. These systems typically leverage individual
 129 LLM strengths by assigning specific roles to each agent, al-
 130 lowing for specialization and improved performance through
 131 in-context learning tailored to specific functions.

132 While a single LLM could theoretically be fine-tuned for
 133 complex planning tasks, this faces key challenges: insuffi-
 134 cient task-specific training data and limited generalization be-
 135 yond the training distribution. Even approaches using com-
 136 plex in-context learning like chain-of-thought prompting can

be unreliable for multi-task scenarios. To address these limi-
 tations, we propose decomposing complex tasks into special-
 ized subtasks handled by different LLM agents working in
 concert. This architecture leverages LLMs' natural strengths
 in language understanding and translation while mitigating
 their weaknesses in complex reasoning tasks.

143 2.2 Planning Formalism

144 A lifted planning task is a tuple $\tau = \langle \mathcal{P}, \mathcal{O}, \mathcal{A}, I, G \rangle$ where 144
 \mathcal{P} is a finite set of first-order predicates and \mathcal{O} a set of ob- 145
 jects. P is a ground predicate or **atom** if all variables have 146
 been replaced by objects. The **goals** G is a set of atoms. A 147
state is a set of atoms; atoms not in the set are assumed to be 148
 false in the state. I is the **initial state**. Each action schema 149
 $A \in \mathcal{A}$ has a list X_A of parameter variables and a **precondi-** 150
tion pre_A , an **add list** add_A and a **delete list** del_A which are 151
 sets of predicates from \mathcal{P} where all variables are replaced by 152
 an element in $X_A \cup \mathcal{O}$. We obtain a (ground) action a from 153
 an action schema A by replacing all variables X_A in pre_A , 154
 add_A , and del_A with an object from \mathcal{O} . A action a is **appli-** 155
icable in a state s if $pre_a \subseteq s$ and $\mathcal{A}(s)$ denotes the set of all 156
 applicable actions in s . Applying action a in state s , results in 157
 the state $s[[a]] = (s \setminus del_a) \cup add_a$. The state resulting from 158
 an iteratively applicable sequence of actions $\pi = \langle a_1, \dots, a_n \rangle$ 159
 is denoted by $s[[\pi]]$. A **plan** is an action sequence π such that 160
 $G \subseteq I[[\pi]]$. A task τ is **solvable** if a plan exists. By $\Pi(\tau)$ we 161
 denote the set of all plans for task τ and by $\tau(G')$ we denote 162
 the task $\tau' = \langle \mathcal{P}, \mathcal{O}, \mathcal{A}, I, G' \rangle$ 163

164 In the following we consider a setting similar to over- 164
 subscription planning [Smith, 2004; Domshlak and Mirkis, 165
 2015], where not all the goals can be satisfied, due for exam- 166
 ple to insufficient resources. But instead of finding a subset 167
 of goals that maximize a utility, we are interested in conflicts 168
 between a set of reference goals G^{ref} . For a task τ these con- 169
 flicts and possible resolutions are given by **minimal unsolv-** 170
able subsets (MUS) [Eifler *et al.*, 2020a] and **minimal cor-** 171
rection sets (MCS) respectively. A set of goals $C \subseteq G^{\text{ref}}$ 172
 is a MUS if $\tau(C)$ is unsolvable but for all $G \subset C$, $\tau(G)$ is 173
 solvable. A set of goals $R \subseteq G^{\text{ref}}$ is a MCS if $\tau(G^{\text{ref}} \setminus R)$ 174
 is solvable but for all $G \subset R$, $\tau(G^{\text{ref}} \setminus G)$ is unsolvable. By 175
 $\mathcal{G}^{\text{MUS}}(\tau, G^{\text{ref}})$ and $\mathcal{G}^{\text{MCS}}(\tau, G^{\text{ref}})$ we denote the set of all MUS 176
 and MCS for task τ with respect to the reference goals G^{ref} . 177
 Both sets can be exponentially large. The following relation 178
 holds between MUS and MCS over the same set of goals G : 179
 $\mathcal{G}^{\text{MCS}}(\tau, G) = \text{HIT}(\mathcal{G}^{\text{MUS}}(\tau, G))$, where $\text{HIT}(\mathcal{S})$ is the set 180
 of all minimal hitting sets of the sets in \mathcal{S} . For algorithms 181
 to compute $\mathcal{G}^{\text{MUS}}(\tau, G)$ we refer to [Eifler *et al.*, 2020a; 182
 Eifler *et al.*, 2020b]. 183

184 Temporal Goals

185 By using a compilation approach [Edelkamp, 2006; Baier and 185
 McIlraith, 2006; De Giacomo *et al.*, 2014] for temporal prop- 186
 erties defined in finite linear temporal logic (LTL_f), it is pos- 187
 sible to reason about conflicts not only over goal facts, but 188
 also over temporal properties of plans [Eifler *et al.*, 2020b]. 189
 In the following we assume that all goals are defined in LTL_f 190
 over atoms. For the syntax and semantic of LTL_f we refer to 191
 [De Giacomo *et al.*, 2014]. With $\mathcal{L}(\tau)$ we denote the set of 192
 all well-formed LTL_f formulas of task τ . 193

194 LTL_f formulas are interpreted over a finite sequence of
 195 states σ . An action sequence π satisfies LTL_f formula ϕ in
 196 state s , iff ϕ holds of the state sequence $\sigma(\pi, s)$ that results
 197 from executing π in s , i.e. iff $\sigma(\pi, s) \models \phi$. A task $\tau(G)$ with
 198 temporal goals G is solvable if there exists an action sequence
 199 π such that for all $\phi \in G$: $\sigma(\pi, I) \models \phi$.

200 3 Iterative Planning with Explanations

201 Iterative planning is based on the idea formalized by [Smith,
 202 2012] that it is often not purposeful to compute just one plan.
 203 Given conflicting goals and users who have not yet fully
 204 formed their preferences, an iterative exploration of possi-
 205 ble plans is more suitable. In this context explanations are
 206 crucial, especially those that help the user to understand the
 207 dependencies between the goals and their preferences. We
 208 define a generic framework for iterative planning with expla-
 209 nations, focussing on the points of interaction with users.

210 The iterative process of determining a final plan is divided
 211 into individual steps. Each step δ_i represents one snapshot of
 212 the user’s exploration of the plan space, defined by a set of
 213 reference goals G_i^{ref} , a set of enforced goals $G_i^{\text{enf}} \subseteq G_i^{\text{ref}}$, and
 214 a sample plan π_i satisfying G_i^{enf} :

215 **Definition 3.1** (Iteration Step). Given a planning task τ , an
 216 **iteration step** is a tuple $\delta = \langle G_i^{\text{ref}}, G_i^{\text{enf}}, \pi \rangle$, where G_i^{ref} is a set
 217 of (temporal) goals for τ , $G_i^{\text{enf}} \subseteq G_i^{\text{ref}}$, and $\pi \in \Pi(\tau(G_i^{\text{enf}}))$ if
 218 $\tau(G_i^{\text{enf}})$ is solvable and $\pi = \epsilon$ otherwise.

219 Based on the sample plan π_i , the user defines the set of
 220 reference goals G_{i+1}^{ref} for the next step. They represent the
 221 goals and preferences that the user is interested in. In addition,
 222 a subset of goals $G_{i+1}^{\text{enf}} \subseteq G_{i+1}^{\text{ref}}$ is selected which must
 223 be satisfied by the sample plan π_{i+1} in the next iteration.

224 The reference goals must be defined in a language the plan-
 225 ner understands; here we are using LTL_f. However, this lan-
 226 guage is not suitable for a lay person as an input language.
 227 Thus, we require a goal translator.

228 **Definition 3.2** (Goal Translator). Given a task τ with well-
 229 formed temporal goals $\mathcal{L}(\tau)$, a **goal translator** is a function
 230 $T_G : NL \mapsto \mathcal{L}(\tau) \cup \epsilon$, that maps a natural language input to
 231 a goal $\phi \in \mathcal{L}(\tau)$ and to ϵ if the natural language description
 232 does not describe a goal represented by any formula in $\mathcal{L}(\tau)$.

233 A more advanced goal translator could include a feedback
 234 loop with the user to recover from misunderstandings or alert
 235 the user that a very similar goal is already considered. For
 236 such features, the translator needs to depend on the *interac-*
 237 *tion context* reflecting the previous interaction with the trans-
 238 lator:

239 **Definition 3.3** (Interaction Context). The **interaction con-**
 240 **text** for translator T is a sequence $\mathcal{C}(T) = c_0 c_1 \dots c_n$ of
 241 interactions $c = \langle \delta, IN, T(IN) \rangle$, where each interaction c is
 242 associated with an iteration step δ and contains the translator
 243 input IN and the translation result $T(IN)$.

244 Explanations can be provided, to facilitate the decision of
 245 which goals to enforce in the next iteration. Those expla-
 246 nations can address different objectives such as clarifying the
 247 model [Sreedharan *et al.*, 2021], identifying the trade-offs be-
 248 tween plan quality measures [Krarup *et al.*, 2024], or provid-
 249 ing a better understanding of the dependencies between the

goals [Eifler *et al.*, 2020a]. Explanations are provided as an- 250
 swers to specific user questions. Depending on whether the 251
 enforced goals G_i^{enf} are satisfiable and therefore a sample plan 252
 π_i exists or whether the enforced goals are unsolvable, the 253
 user’s questions will vary. In the former case, the question 254
 may relate to how the sample plan solves the task whereas in 255
 the latter, the user is more interested in why G_i^{enf} cannot be 256
 satisfied. Formally we define a question as follows. 257

258 **Definition 3.4** (Question). Given a task τ and a set of ques-
 259 tion types M_Q , a **question** is a tuple $Q = \langle \mu_Q, args \rangle$ where
 260 $\mu_Q \in M_Q$ is the question type and $args \subseteq \wp(\mathcal{L}(\tau))$ (the
 261 powerset of $\mathcal{L}(\tau)$) are the question arguments.

262 The question types M_Q depend on the explanation frame-
 263 work. In Section 4 we introduce the types our explanation
 264 framework supports. An example is $\mu_Q = \text{S-why-not}$, i. e.
 265 “Why is g not satisfied by plan π ?”, which requires one argu-
 266 ment. Roughly, each question word in natural language, e. g.
 267 *who, what, why, how* maps to a different type of question.

268 This is the second point of interaction. Again, we require
 269 a translator to allow the user to ask their question in natural
 270 language, which is then translated into a formal question.

271 **Definition 3.5** (Question Translator). Given a task τ , and a
 272 set of question types M_Q , a **question translator** is a function
 273 $T_Q : NL \mapsto (M_Q \times \wp(\mathcal{L}(\tau))) \cup \epsilon$, that maps a natural language
 274 expression to a question type and its arguments and to ϵ if no
 275 matching question type exists.

276 Question translation can be divided into two steps. First
 277 the question type is identified. Then, the question topic, i. e.
 278 the arguments, are translated. This last step can be performed
 279 by the goal translator, as the translation tasks are the same.

280 A context-dependent question translator with access to pre-
 281 viously asked questions enables follow-up questions with im-
 282 plicit references. For example, for the questions “Why can I
 283 not visit Alice?” and “Can you enforce it?”, the visited lo-
 284 cation or even the entire question argument depend on the
 285 context.

286 The explanation framework computes a set of formal ex-
 287 planations $\mathcal{E}(Q)$ based on the translated question $Q =$
 288 $\langle \mu_Q, args \rangle$ and all additional required data, for example the
 289 planning task τ and the current iteration step δ_i . We do not
 290 place any restrictions on the exact format or language of an
 291 explanation $E \in \mathcal{E}(Q)$, and will simply refer to the language
 292 as $\mathcal{L}_{\mathcal{E}}$. However, we assume that each explanation $E \in \mathcal{E}(Q)$
 293 is sufficient in itself to answer the question. $\mathcal{E}(Q)$ is regarded
 294 as a selection of possible explanations.

295 These explanations must be communicated to the user, and
 296 thus we again need a translator function.

297 **Definition 3.6** (Explanation Translator). An **explanation**
 298 **translator** is a function $T_{\mathcal{E}} : \wp(\mathcal{L}_{\mathcal{E}}) \mapsto NL$, that maps a
 299 set of formal explanations to a natural language explanation.

300 This is a simple version of an explanation translator that
 301 does not offer a user or context-dependent translation. User-
 302 dependent translation is important in order to customize the
 303 vocabulary to the user and to provide an answer with the ex-
 304 pected level of detail. Incorporating the interaction context
 305 can have several benefits, from the possibility to ask follow-
 306 up questions, which naturally increases interactivity, to the

307 inclusion of previous interactions into the selection and sum-
308 marization of the explanations.

309 **Selecting Explanations** The explanation framework pro-
310 vides a selection of explanations, each of which answers the
311 question Q on its own. However, some causes or properties
312 may not be as relevant as others or may not fit into the in-
313 teraction context. In line with the insight from social science
314 [Miller, 2019] that humans select small relevant explanations
315 given the context, it can therefore be advantageous to use only
316 a subset of explanations.

317 Often the size of the explanations is chosen as selection
318 criteria [Chakraborti *et al.*, 2017]. [Junker, 2004] selects *pre-*
319 *ferred* explanations based on an upstream or interleaved pro-
320 cess to collect a preference ranking over the goals.

321 **Summarizing Explanations** It may be preferable to con-
322 vey multiple explanations $\hat{\mathcal{E}} \subseteq \mathcal{E}(Q)$, in cases where the ef-
323 fect has multiple sufficient causes neither of which is neces-
324 sary. But since explanations may be at different levels and
325 expressing their logical relation results in an overly lengthy
326 response, it is often simpler to provide a summary of the
327 explanations $\hat{\mathcal{E}}$. In line with how humans give explanations
328 [Miller, 2019], such a summary should convey only relevant
329 information as effectively as possible in the given context. As
330 for the explanation selection, the user questions can provide
331 an indication of the level of detail expected. Also, the in-
332 teraction context can be used to enable or facilitate the sum-
333 marization by referring to previously addressed questions or
334 explanations.

335 Selecting the correct abstraction level for the explanations
336 [Sreedharan *et al.*, 2019], or only communicating via a pre-
337 defined vocabulary [Vasileiou and Yeoh, 2023] also leads to a
338 summarization in the sense of leaving out details that are not
339 required or known.

340 4 Goal Conflict Explanations

341 In the following, we extend the explanation framework by
342 [Eifler *et al.*, 2020a; Eifler *et al.*, 2020b] to address a larger
343 number of question types, allowing a richer user interaction.
344 The questions types have been collected during an interview
345 conducted at a manufacturing company¹ with the end users
346 of a future explainable planning system. We call this exten-
347 sion *Explanation Framework with Conflicts and Corrections*
348 (EF_{CC}), in line with the fact that all answers are based on
349 either the minimal conflicts (MUS) or minimal corrections
350 (MCS) of the reference goals G^{ref} and the arguments of the
351 question. In the following definitions we assume the user asks
352 questions only about goals already included in G^{ref} . If this is
353 not the case then they can be simply added to G^{ref} before
354 calling the MUS or MCS computation. We list all supported
355 question types and the corresponding formal explanations. To
356 clarify the meaning of the question and answer we include
357 one possible natural language version of both. As an input
358 we require the task τ , the iteration step $\delta = \langle G^{\text{ref}}, G^{\text{enf}}, \pi \rangle$
359 and the question $Q = \langle \mu_Q, args \rangle$. The produced explanations
360 $\mathcal{E}(Q)$ are subsets of G^{ref} , which means the formal language
361 of an explanation is $\mathcal{L}_{\mathcal{E}} = \wp(\mathcal{L}(\tau))$.

¹The company is not named to preserve submission anonymity.

If $\tau(G^{\text{enf}})$ is unsolvable we support the following two
question types:

US-why: “Why is the task unsolvable?”:

$$\mathcal{E}(\langle \text{US-why}, \emptyset \rangle) = \mathcal{G}^{\text{MUS}}(\tau, G^{\text{enf}})$$

- “The task is unsolvable because it is not possible to sat-
isfy any of the conflicts in $\mathcal{E}(\langle \text{US-why}, \emptyset \rangle)$.”

US-how: “How can I make the task solvable?”

$$\mathcal{E}(\langle \text{US-how}, \emptyset \rangle) = \mathcal{G}^{\text{MCS}}(\tau, G^{\text{enf}})$$

- “To make the task solvable you have to forego one of the
goal sets in $\mathcal{E}(\langle \text{US-how}, \emptyset \rangle)$.”

If $\tau(G^{\text{enf}})$ is solvable, we support question types referring to
goals not satisfied by the sample plan π . By $G^{\text{true}}(\pi) = \{\phi \in$
 $G^{\text{ref}} \mid \sigma(\pi, I) \models \phi\}$ we denote the goals satisfied by π , and
by $G^{\text{false}}(\pi) = G^{\text{ref}} \setminus G^{\text{true}}(\pi)$ the goals not satisfied by π .
For all the following question types, the arguments must not
be satisfied by the current plan, i.e. $args \subseteq G^{\text{false}}(\pi)$.

Answers to question types $\{S\text{-why-not}, S\text{-what-if},$
 $S\text{-can}\}$ are based on the same information, but are phrased
differently.

$$\mathcal{E}(\langle \mu_Q, args \rangle) = \min_{\subseteq}(\{C \setminus args \mid$$

 $C \in \mathcal{G}^{\text{MUS}}(\tau, G^{\text{ref}}), C \subseteq G^{\text{true}}(\pi) \cup args\})$

where $\min_{\subseteq}(S)$ filters the sets in S that are subset-minimal.

S-why-not: “Why are $args$ not satisfied?”

- $\mathcal{E}(\langle S\text{-why-not}, args \rangle) = \emptyset$: “ $args$ can be satisfied
without foregoing any of the already satisfied goals.”
- $\emptyset \in \mathcal{E}(\langle S\text{-why-not}, args \rangle)$: “The goals in $args$ cannot
be satisfied together.”
- otherwise: “There is a conflict between $args$ and all the
goal subsets in $\mathcal{E}(\langle S\text{-why-not}, args \rangle)$.”

S-what-if: “What happens if we enforce $args$?”

- $\mathcal{E}(\langle S\text{-what-if}, args \rangle) = \emptyset$: “ $args$ can be satisfied
without foregoing any goal satisfied by the plan.”
- $\emptyset \in \mathcal{E}(\langle S\text{-what-if}, args \rangle)$: “Then the problem would
be unsolvable.”
- otherwise: “You could no longer satisfy any of the goal
sets in $\mathcal{E}(\langle S\text{-what-if}, args \rangle)$.”

In the yes/no question type $S\text{-can}$, the intention that none
of the currently satisfied goals should be given up is implicit.
This interpretation is based on the sample responses of the
human planners who did not consider the option of foregoing
any satisfied goals in $G^{\text{true}}(\pi)$.

S-can: “Can $args$ be satisfied?”

- $\mathcal{E}(\langle S\text{-can}, args \rangle) = \emptyset$: “ $args$ can be satisfied.”
- otherwise: “It is not possible.”

S-how: “How can $args$ be satisfied?”

$$\mathcal{E}(\langle S\text{-how}, args \rangle) = \min_{\subseteq}(\{C \setminus G^{\text{false}}(\pi) \mid$$

 $C \in \mathcal{G}^{\text{MCS}}(\tau, G^{\text{ref}}), C \cap args = \emptyset\})$

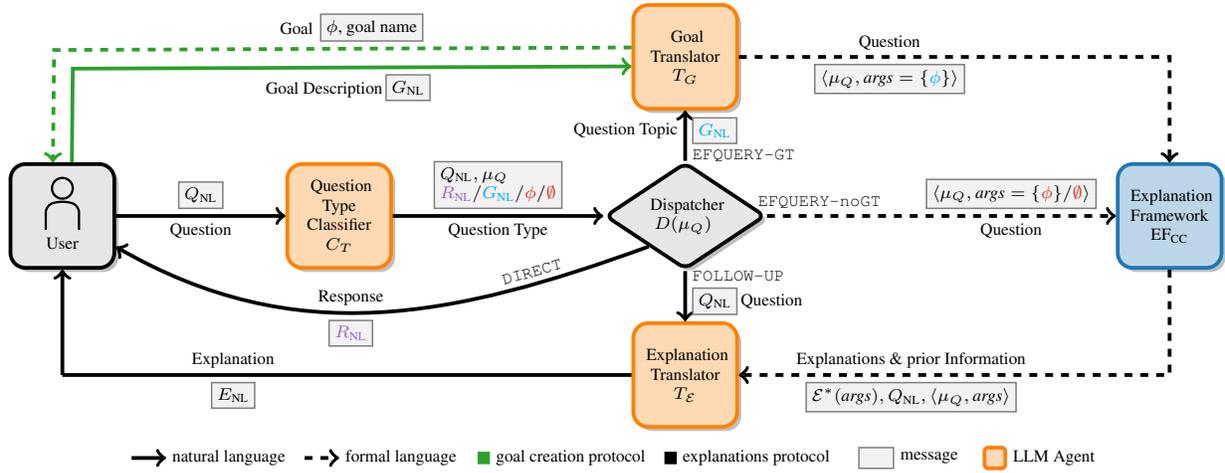


Figure 2: Communication protocol between the user, the translators and the explanation framework. In addition to the information provided in the messages, each agent has access to the planning task and the iteration step. For a goal translation (green, top left) the user directly communicates with T_G . If the user asks a question (black) the dispatcher chooses one of four routing options (DIRECT, FOLLOW-UP, EFQUERY-noGT and EFQUERY-GT) depending on the question type identified by C_T .

- 407 • $\emptyset \in \mathcal{E}(\langle S\text{-how}, args \rangle)$: “args can be satisfied without
408 foregoing any goals satisfied by the plan.
- 409 • $\mathcal{E}(\langle S\text{-how}, args \rangle) = \emptyset$: “It is not possible.”
- 410 • otherwise: “You have to forego one of the goal sets in
411 $\mathcal{E}(\langle S\text{-how}, args \rangle)$.”

412 The number of MUS and MCS can be exponential in the number
413 of goals, but for all question types M_Q one goal subset
414 $E \in \mathcal{E}(Q)$ is sufficient to answer the question. However,
415 some conflicts may be more relevant, or it may be easier to
416 forego some corrections. Therefore, a selection of $\hat{\mathcal{E}} \subseteq \mathcal{E}(Q)$
417 and a summarization of the explanations $\hat{\mathcal{E}}$ is desirable.

418 5 Multi-Agent LLMs

419 Next, we describe the implementation of our interactive
420 framework using LLM based translators to communicate with
421 EF_{CC}. Each translator, as well as the planner and EF_{CC} (which
422 we consider black boxes that generate plans and explanations
423 when given a correct input message), are *agents* that commu-
424 nicate with each other. We first present the communication
425 protocol between these agents, and then discuss the specifics
426 of each translator agent.

427 5.1 Communication Protocol

428 Figure 2 shows the communication protocol between the
429 user, translators and the explanation framework. To trans-
430 late a natural language goal description G_{NL} into an LTL_f
431 expression ϕ , the user directly communicates with the goal
432 translator T_G .

433 When a user asks a question Q_{NL} , different routes are
434 possible. The *Question Translator* T_Q is composed of two
435 agents: the *Question Type Classifier* C_T and the *Goal Trans-*
436 *lator* T_G . A *Dispatcher* routes the response of C_T depend-
437 ing on the identified question type μ_Q .

438 $\mu_Q = \text{DIRECT}$: If C_T is capable of answering the question,
439 then the answer is given directly to the user. This includes for

440 example questions like “Which questions can I ask?” but also
441 incomprehensible messages or unsupported questions. More
442 examples are given in the technical appendix.

443 $\mu_Q = \text{FOLLOW-UP}$: If Q_{NL} is identified as a follow-up ques-
444 tion, not requiring new input from EF_{CC}, then Q_{NL} is directly
445 sent to the *Explanation Translator* T_E .

446 $\mu_Q \in M_Q$: If Q_{NL} is classified as one of the question
447 types in M_Q , then the question topic is forwarded to T_G
448 and the combined result is sent to the explanation framework
449 (EFQUERY-GT). If the question type does not require an argu-
450 ment or when the question topic is about a goal already
451 used in a previous iteration step (i.e., in $G_{all}^{ref} = \bigcup G_i^{ref}$) T_G
452 is bypassed and the (already known) LTL_f formula correspond-
453 ing to the goal is used (EFQUERY-noGT). The resulting formal
454 explanation is then translated by T_E .

455 5.2 Translator Implementations

456 All agents are instantiated from the same base model *GPT-*
457 *4o-mini*. Their contexts, the input and output histories, are
458 maintained separately from each other. Each LLM agent re-
459 ceives input in the form of a structured string consisting of
460 the prompts described in the technical appendix and the input
461 components listed in Figure 2.

462 Next we address the tasks of the individual agents and the
463 specific details of their implementation.

464 **Question Type Classifier** The *Question Type Classifier* C_T
465 is the first step in question processing. The tasks of this agent
466 extend beyond the identification of the question type μ_Q . In
467 case of DIRECT it directly generates the response R_{NL} . For
468 EFQUERY-GT and EFQUERY-noGT, it extracts a goal de-
469 scription G_{NL} and identifies if G_{NL} matches with a goal de-
470 scription $G_{NL} \in G_{all}^{ref}$. It can also identify FOLLOW-UP ques-
471 tions to route them directly to T_E .

472 **Goal Translator** The goal translator T_G produces both the
473 LTL_f formula and a natural language description of the goal
474 (used in UI and communication between LLM agents).

475 The task of translating natural language to LTL or LTL_f
476 has been explored in different fields with different approaches
477 [Brunello *et al.*, 2019]. More recently, approaches based on
478 LLMs that use prompting have been used to implement tools
479 such as NL2LTL [Fuggitti and Chakraborti, 2023] a template-
480 based classifier, Lang2LTL [Liu *et al.*, 2022] which works
481 without templates and only provides the available literals to
482 the LLM, and nl2spec [Cosler *et al.*, 2023] which ad-
483 dresses sub-formulas iteratively to counteract ambiguities.

484 Our evaluation focuses on the question classification and
485 the explanation translation. Thus, we opted for simple base
486 implementation, similar to what [Liu *et al.*, 2022] call *End-*
487 *to-End Approach*. The LLM is provided with the predicates
488 and objects that can compose the literals and a few samples.

489 **Explanation Translator** An advantage of an LLM-based
490 translator is its ability to provide user-dependent translation,
491 select explanations, and summarize them. Instead of just pro-
492 viding explanations, we provide both conflicts and correc-
493 tions $\mathcal{E}^*(Q)$. This is motivated by the anticipation of follow-
494 up questions (e.g. a “how” question after a “why” question).
495 It additionally receives the user’s question Q_{NL} , and the for-
496 mal translation Q . μ_Q indicates how to interpret $\mathcal{E}^*(Q)$,
497 while Q_{NL} helps to generate a natural response. We expect
498 the explanation translator to leverage prompt examples to un-
499 derstand how to select or summarize explanations. It is thus a
500 domain-specific design choice to provide examples of the ex-
501 pected selection/summarization strategy. With FOLLOW-UP
502 questions, users can actively request a summary or selection.
503 We leave the evaluation of this feature and the extraction of
504 user preferences from selection requests as future work.

505 Context-Dependent Translators

506 A translator solely based on the inputs specified in Figure 2,
507 can only provide context-independent translations and is not
508 capable of addressing follow-up question. To enable context-
509 dependent translators we leverage the context memory of
510 LLMs: they retain their context across interactions, so they
511 can build on and reference previous interactions. This in-
512 herently provides all the benefits listed in Section 3. While
513 LLM contexts can theoretically retain every previous itera-
514 tions with a user, we implemented separate contexts for each
515 iteration step. This allows to maintain clarity by avoiding
516 mixing information from different iteration steps. Examples
517 of context-dependant interactions, selections or summariza-
518 tion of explanations are given in the technical appendix.

519 6 Evaluation

520 We developed a web-based platform extending [Eifler *et al.*,
521 2022] with our LLM-based framework and an improved UI.
522 The source code will be published upon acceptance. We then
523 used this tool to conduct a pilot user study aiming at evaluat-
524 ing the effectiveness of some of our LLM agents and the im-
525 pact of our LLM-based explanation interface on users’ ability
526 to solve a planning task. The material presented in this sec-
527 tion is complemented in appendix.

528 6.1 User Study Design

529 Similarly to [Eifler *et al.*, 2022], we used a parent’s afternoon
530 planning task as the scenario. The participants were tasked

531 with the planning of afternoon activities for a family, not all
532 of which can be carried out due to time constraints. As a
533 proxy for user preference, each goal has an associated util-
534 ity. The objective of the participants was to select a subset of
535 the goals to maximize the total utility of the plan. The goal
536 utility is unknown to all other agents (planner, explainer and
537 LLMs). Consequently, these agents are capable of providing
538 information regarding the solvability of a set of goals, and
539 the conflicts and corrections, but are unable to provide direct
540 assistance in maximizing utility.

541 The study was divided into three parts. First, participants
542 were introduced to the tool through an introduction task (6
543 goals). Second, a more complex task (15 goals with 35 con-
544 flicts and 19 corrections) was used to compare the effect of
545 the different interfaces used by G^{TPL} and G^{LLM} on maximizing
546 the utility within 20 min. Finally, participants completed a
547 post-experiment questionnaire.

548 In order to evaluate the effectiveness of LLM agents in an
549 interactive explanation framework, we divided the 70 partic-
550 ipants into two groups. Our control group G^{TPL} included 34
551 participants utilizing a **template-based** interface, which al-
552 lowed them to select from the predefined set of questions in-
553 troduced in Section 4 and receive answers based on simple
554 natural language templates. The second group G^{LLM} included
555 36 participants given the **LLM-based** interface, with which
556 they asked freely formulated natural language questions that
557 were processed by LLM-based translators as described in
558 Section 5. We did not include a group without explanations
559 as an earlier study [Eifler *et al.*, 2022] had already examined
560 that baseline and demonstrated that our control condition was
561 more effective.

562 6.2 Evaluation of LLM Agents

563 We first evaluated the accuracy of the question type clas-
564 sifier C_T on the complete set of 134 questions collected
565 from G^{LLM} . The classification of both the question type and
566 the question argument was correct in 89.47% of cases. All
567 observed failures were related to routing decisions, where
568 the question translator T_Q unnecessarily routed questions to
569 EFCC by predicting $\mu_Q \in M_Q$ instead of responding directly
570 (DIRECT) or forwarding them to the explanation translator
571 T_E (FOLLOW-UP). Importantly, these classification errors re-
572 sulted in suboptimal routing – with no delay since the expla-
573 nations were already computed – but T_E was still able to pro-
574 vide appropriate responses to the users’ questions.

575 To ensure a controlled user study environment with reason-
576 able latency for explanations, we decided to provide a fixed
577 set of reference goals. This results in all users having the
578 same optimization task and allows to precompute all explana-
579 tions. This means that users did not create new goals, which
580 leaves the evaluation of the goal translation for future work.
581 As for T_E , the main challenge is to measure the correctness of
582 a summarized explanation which is out of scope of this paper.

583 6.3 User Groups Comparison

584 We compared the two groups using the questionnaire and a
585 number of metrics, including the time spent on the task, the
586 maximum utility reached, and the number of questions asked.

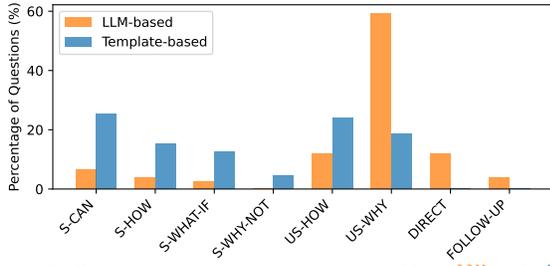


Figure 3: Comparison of question types used by G^{LLM} and G^{TPL} .

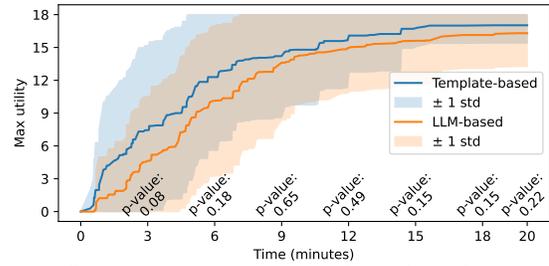


Figure 4: Comparison of maximum utility achieved over time between G^{LLM} and G^{TPL} . Group means are computed at each time step.

587 Users took on average 12.0 mins \pm 7.0 in G^{LLM} and 12.3
 588 mins \pm 6.3 in G^{TPL} to complete the evaluation task. 22/36
 589 (61.1%) users in G^{LLM} and 23/34 (67.6%) in G^{TPL} reached
 590 the maximum utility of 18.

591 G^{TPL} asked on average significantly more questions (10.5
 592 \pm 13.3) vs. (2.5 \pm 2.6) for G^{LLM} ($p < 0.001$, t-test). G^{TPL} also
 593 asked more diverse questions, as shown in Figure 3. Although
 594 G^{LLM} had complete freedom in the questions they could ask,
 595 they mainly asked US-why questions. The analysis of partic-
 596 ipants’ feedback on their strategies also shows more diversity
 597 in the G^{TPL} group. These points suggest that the G^{LLM} users
 598 did not take full advantage of the conversational capabilities
 599 of the LLM interface.

600 When comparing the evolution of the plan utility over time
 601 shown in Figure 4, we observe that G^{TPL} has a higher mean
 602 value than G^{LLM} for the whole duration of the task. Nonethe-
 603 less, the difference is not statistically significant. The utility
 604 over iteration steps lead to similar conclusions, and shows a
 605 convergence significantly quicker toward the maximum util-
 606 ity for G^{TPL} . The questionnaire results only shows significant
 607 difference on the question “Did questions help to improve a
 608 plan?” which received higher scores from G^{TPL} . This cor-
 609 roborates the fact that G^{LLM} users mostly used the explanation
 610 interface during unsolvable iteration steps.

6.4 Limitations and Discussion

612 **Instance and Task** While this user study provides some
 613 evidence of better performance of G^{TPL} over G^{LLM} , it is im-
 614 portant to note that these results were obtained on a single
 615 planning task instance of moderate difficulty. Although this
 616 setup was chosen to align with a previous user study in which
 617 template-based goal-conflict explanations were found useful
 618 [Eifler *et al.*, 2022], it remains to investigate how the pro-
 619 posed architecture would help on a task for which template-
 620 based approaches show limitations — particularly problems
 621 where summarization is a *necessity*. Perhaps a more extensive
 622 training session, highlighting the LLM’s capabilities, could
 623 have helped users to take fuller advantage of the interface.
 624 Another limitation of our study is that it was conducted with
 625 lay users. Whether similar results hold for the target users of
 626 such tools (domain experts) remains open. One can hypoth-
 627 esize that domain experts might ask more diverse questions
 628 than lay persons.

629 **Using LLMs has a (time and effort) cost** The free-text
 630 feedback obtained through the questionnaire provides in-
 631 sights into why G^{LLM} performed slightly worse than G^{TPL} .
 632 Common concerns included system latency (response times

633 averaged \sim 5 secs, occasionally exceeding 10 secs) and the
 634 additional effort of formulating and typing questions. In par-
 635 ticular, G^{LLM} needed a much deeper understanding of which
 636 questions to ask, whereas G^{TPL} readily had access to the rel-
 637 evant questions. These factors led G^{LLM} users to ask fewer
 638 questions and explore less thoroughly the range of possi-
 639 ble interactions. Additionally, receiving an unconvincing re-
 640 sponse early in the task could diminish trust in the system –
 641 an effect that particularly impacts G^{LLM} users since the expla-
 642 nation interface requires more effort and patience.

7 Conclusion and Future Work

643 **Contributions** In this paper we presented a high-level ar-
 644 chitecture for interactive explanations in the context of iter-
 645 ative planning. We instantiated this framework with goal-
 646 conflict explanations and provided an easy-to-use implemen-
 647 tation via a web platform. Finally, we conducted a user study
 648 to evaluate the effect of the LLM-based explanation interface
 649 over a more conventional template-based interface.

651 **Findings** We found that there is some potential for LLM-
 652 based translators but they did not show objective improve-
 653 ments in helping lay users to solve a moderately difficult plan-
 654 ning task. Users of the LLM interface asked fewer and less
 655 diverse questions than the users of the template-based inter-
 656 face. While this could be due to the additional effort and
 657 patience required to use the LLM interface, it also indicates
 658 that users of the template-based interface had a better under-
 659 standing of the capabilities of the system. The fact that we
 660 did not observe this with expert users that tested this tool and
 661 successfully utilized the capabilities of the LLM suggests that
 662 there is a learning curve that some user study participants did
 663 not overcome in time, and that LLM-based interfaces could
 664 be helpful in other explainable planning contexts.

665 **Future Work** Future work should evaluate LLM-translated
 666 explanations on tasks on which template-based translations
 667 show limitations. Summarization capabilities of LLMs in this
 668 context should be assessed carefully.

669 We also plan on investigating the extension of this ap-
 670 proach to other explanations frameworks. This includes
 671 MUS and MCS based approaches used in Constraint Pro-
 672 gramming [Povéda *et al.*, 2024; Gamba *et al.*, 2023] but
 673 also other approaches used in planning [Krarup *et al.*, 2021;
 674 Sreedharan *et al.*, 2021]. In addition to the classification of
 675 the question type, this adds the challenge of deciding which
 676 explanation approach is most suitable to answer the question.

677 Ethical Statement

678 We confirm that our user study was conducted under full compliance with the requirements of our local ethics committee [name withheld for anonymity]. All participants were volunteers who gave informed consent, and the study was designed more like a playful, game-like engagement than a constraining activity, ensuring no potential harm or emotional discomfort. Furthermore, no sensitive elements were introduced that could adversely affect the participants' well-being.

686 While large language models (LLMs) in our framework may occasionally generate misleading or irrelevant statements (i.e., "hallucinations"), the core planning tasks are handled by a dedicated planner, eliminating the risk of producing incorrect plans. At worst, iterative planning might be rendered less effective if LLM responses deviate from the actual planning context. Finally, our emphasis on explainability furthers transparency, aligning our work with the broader objective of making AI-driven systems more interpretable and trustworthy.

696 References

697 [Achinstejn, 1980] Peter Achinstejn. *The Nature of Explanation*. Oxford University Press, 1980.

699 [Baier and McIlraith, 2006] Jorge A. Baier and Sheila A. McIlraith. Planning with first-order temporally extended goals using heuristic search. In *Proc. AAAI*, 2006.

702 [Bommasani *et al.*, 2021] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

708 [Bromberger, 1962] Sylvain Bromberger. An approach to explanation. In R. Butler, editor, *Analytical Philosophy*. Oxford University Press, 1962.

711 [Brown *et al.*, 2020] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Proc. NeurIPS*, 2020.

716 [Brunello *et al.*, 2019] Andrea Brunello, Angelo Montanari, and Mark Reynolds. Synthesis of ltl formulas from natural language texts: State of the art and research directions. In *Proc. TIME*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.

721 [Chakraborti *et al.*, 2017] Tathagata Chakraborti, Sarath Sreedharan, Yu Zhang, and Subbarao Kambhampati. Plan explanations as model reconciliation: Moving beyond explanation as soliloquy. In *IJCAI*, 2017.

725 [Cosler *et al.*, 2023] Matthias Cosler, Christopher Hahn, Daniel Mendoza, Frederik Schmitt, and Caroline Trippel. nl2spec: interactively translating unstructured natural language to temporal logics with large language models. In *Proc. CAV*, 2023.

[Dazeley *et al.*, 2021] Richard Dazeley, Peter Vamplew, Cameron Foale, Charlotte Young, Sunil Aryal, and Francisco Cruz. Levels of explainable artificial intelligence for human-aligned conversational explanations. *AIJ*, 2021.

[De Giacomo *et al.*, 2014] Giuseppe De Giacomo, Riccardo De Masellis, and Marco Montali. Reasoning on LTL on finite traces: Insensitivity to infiniteness. In *AAAI*, 2014.

[Domshlak and Mirkis, 2015] Carmel Domshlak and Vitaly Mirkis. Deterministic oversubscription planning as heuristic search: Abstractions and reformulations. *JAIR*, 2015.

[Edelkamp, 2006] Stefan Edelkamp. On the compilation of plan constraints and preferences. In *ICAPS*, 2006.

[Eifler *et al.*, 2020a] Rebecca Eifler, Michael Cashmore, Jörg Hoffmann, Daniele Magazzeni, and Marcel Steinmetz. A new approach to plan-space explanation: Analyzing plan-property dependencies in oversubscription planning. In *AAAI*, 2020.

[Eifler *et al.*, 2020b] Rebecca Eifler, Marcel Steinmetz, Alvaro Torralba, and Jörg Hoffmann. Plan-space explanation via plan-property dependencies: Faster algorithms & more powerful properties. In *IJCAI*, 2020.

[Eifler *et al.*, 2022] Rebecca Eifler, Martim Brandao, Amanda Coles, Jeremy Frank, and Jörg Hoffmann. Evaluating plan-property dependencies: A web-based platform and user study. In *ICAPS*, 2022.

[Feldhus *et al.*, 2023] Nils Feldhus, Qianli Wang, Tatiana Anikina, Sahil Chopra, Cennet Oguz, and Sebastian Möller. Interrolang: Exploring nlp models and datasets through dialogue-based explanations. In *Proc. EMNLP*, 2023.

[Fuggitti and Chakraborti, 2023] Francesco Fuggitti and Tathagata Chakraborti. NL2LTL – a python package for converting natural language (NL) instructions to linear temporal logic (LTL) formulas. In *ICAPS*, 2023.

[Gamba *et al.*, 2023] Emilio Gamba, Bart Bogaerts, and Tias Guns. Efficiently explaining cpsps with unsatisfiable subset optimization. *JAIR*, 2023.

[Guo *et al.*, 2024] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*, 2024.

[Junker, 2004] Ulrich Junker. Preferred explanations and relaxations for over-constrained problems. In *Proc. AAAI*, 2004.

[Kambhampati *et al.*, 2024] Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri, Lucas Saldyt, and Anil Murthy. Position: Llms can't plan, but can help planning in llm-modulo frameworks. In *ICML*, 2024.

[Krarup *et al.*, 2021] Benjamin Krarup, Senka Krivic, Daniele Magazzeni, Derek Long, Michael Cashmore, and David E. Smith. Contrastive explanations of plans through model restrictions. *JAIR*, 2021.

- 784 [Krarup *et al.*, 2024] Benjamin Krarup, Amanda Jane Coles,
785 Derek Long, and David E. Smith. Explaining plan quality
786 differences. In *ICAPS*, 2024.
- 787 [Lakkaraju *et al.*, 2022] Himabindu Lakkaraju, Dylan Slack,
788 Yuxin Chen, Chenhao Tan, and Sameer Singh. Rethinking
789 explainability as a dialogue: A practitioner’s perspective.
790 In *NeurIPS Workshop on Human Centered AI*, 2022.
- 791 [Liao *et al.*, 2020] Q Vera Liao, Daniel Gruen, and Sarah
792 Miller. Questioning the ai: informing design practices for
793 explainable ai user experiences. In *Proc. CHI*, 2020.
- 794 [Liu *et al.*, 2022] Jason Xinyu Liu, Ziyi Yang, Benjamin
795 Schornstein, Sam Liang, Ifrah Idrees, Stefanie Tellex, and
796 Ankit Shah. Lang2tl: Translating natural language com-
797 mands to temporal specification with large language mod-
798 els. In *Workshop on Language and Robotics at CoRL*,
799 2022.
- 800 [Miller, 2019] Tim Miller. Explanation in artificial intelli-
801 gence: Insights from the social sciences. *AI*, 2019.
- 802 [Nguyen *et al.*, 2023] Van Bach Nguyen, Jörg Schlötterer,
803 and Christin Seifert. From black boxes to conversations:
804 Incorporating xai in a conversational agent. In *XI*, 2023.
- 805 [Peng *et al.*, 2024] Binghui Peng, Sridhar Narayanan, and
806 Christos Papadimitriou. On limitations of the transformer
807 architecture. *arXiv preprint arXiv:2402.08164*, 2024.
- 808 [Povéda *et al.*, 2024] Guillaume Povéda, Andreas Strahl,
809 Mark Hall, Ryma Boumazouza, Santiago Quintana-
810 Amate, Nahum Alvarez, Ignace Bleukx, Dimos Tsouros,
811 H el ene Verhaeghe, and Tias Guns. Trustworthy and ex-
812 plainable decision-making for workforce allocation. In
813 *Workshop on Progress Towards the Holy Grail at CP*,
814 2024.
- 815 [Shen *et al.*, 2023] Hua Shen, Chieh-Yang Huang, Tong-
816 shuang Wu, and Ting-Hao Kenneth Huang. Convxai: De-
817 livering heterogeneous AI explanations via conversations
818 to support human-ai scientific writing. *CoRR*, 2023.
- 819 [Slack *et al.*, 2023] Dylan Slack, Satyapriya Krishna,
820 Himabindu Lakkaraju, and Sameer Singh. Explaining
821 machine learning models with interactive natural lan-
822 guage conversations using talktomodel. *Nature Machine*
823 *Intelligence*, 2023.
- 824 [Smith, 2004] David E. Smith. Choosing objectives in over-
825 subscription planning. In *ICAPS*, 2004.
- 826 [Smith, 2012] David Smith. Planning as an iterative process.
827 In *AAAI*, 2012.
- 828 [Sreedharan *et al.*, 2019] Sarath Sreedharan, Siddharth Sri-
829 vastava, David E. Smith, and Subbarao Kambhampati.
830 Why can’t you do that hal? explaining unsolvability of
831 planning tasks. In *IJCAI*, 2019.
- 832 [Sreedharan *et al.*, 2021] Sarath Sreedharan, Tathagata
833 Chakraborti, and Subbarao Kambhampati. Foundations of
834 explanations as model reconciliation. *AIJ*, 2021.
- 835 [State *et al.*, 2023] Laura State, Salvatore Ruggieri, and
836 Franco Turini. Reason to explain: Interactive contrastive
837 explanations (reasonx). In *XAI*, 2023.
- [Touvron *et al.*, 2023] Hugo Touvron, Thibaut Lavril, Gau- 838
tier Izacard, Xavier Martinet, Marie-Anne Lachaux, Tim- 839
oth e Lacroix, Baptiste Rozi re, Naman Goyal, Eric Ham- 840
bro, Faisal Azhar, et al. Llama: Open and efficient founda- 841
tion language models. *arXiv preprint arXiv:2302.13971*, 842
2023. 843
- [Valmeekam *et al.*, 2023] Karthik Valmeekam, Matthew 844
Marquez, Sarath Sreedharan, and Subbarao Kambham- 845
pati. On the planning abilities of large language models-a 846
critical investigation. *NeurIPS*, 2023. 847
- [Vasileiou and Yeoh, 2023] Stylianos Loukas Vasileiou and 848
William Yeoh. PLEASE: generating personalized expla- 849
nations in human-aware planning. In *Proc. ECAI*, 2023. 850
- [Vaswani *et al.*, 2017] Ashish Vaswani, Noam Shazeer, Niki 851
Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, 852
Lukasz Kaiser, and Illia Polosukhin. Attention is all you 853
need. In *NeurIPS*, 2017. 854
- [Zhang *et al.*, 2024] Tong Zhang, X Jessie Yang, and Boyang 855
Li. May i ask a follow-up question? understanding the 856
benefits of conversations in neural network explainabil- 857
ity. *International Journal of Human-Computer Interac-* 858
tion, 2024. 859