

# ProcessChat: A Business Process Grounded Dialogue Dataset

## Anonymous ACL submission

### Abstract

Business processes are designed to streamline and optimize work within an organization and are often defined and documented by domain experts or process analysts using formal specifications. However, these specifications may be complex for the users executing the tasks of the process. For example, a recruitment process designed by a domain expert is used by many actors in the organization, who may not be skilled in understanding the formal notations that specify the process. With recent advancements in large language models, there has been increasing interest in enabling users to ask questions in natural language and receive relevant responses that are specific to the user’s context and process knowledge. We propose a dialog dataset grounded in domain-specific process knowledge, which it is supposed to follow during the conversation. The dataset consists of 316 dialogs grounded on 73 different process model specifications. We also present a baseline model, which is trained on the proposed dataset. Our experiments find that the model can do zero-shot transfer to unseen processes, and sets a strong baseline for future research.

## 1 Introduction

A commonly used standard for process modeling is the Business Process Model and Notation (BPMN), overseen by the Object Management Group (OMG) which provides a rich set of notations (OMG, 2011) to represent the process operations. While process models (or specifications) are useful artifacts to represent the operations, one of the challenges in using them is that they are not intuitive to business professionals, who conduct various tasks of the process (van der Aa et al., 2015). For example, a Human Resource (HR) process for a leave application can be defined by a domain expert using BPMN specification. Figure 1 presents a simple leave application process specified using BPMN. The specification comprises: i) the activities (rect-

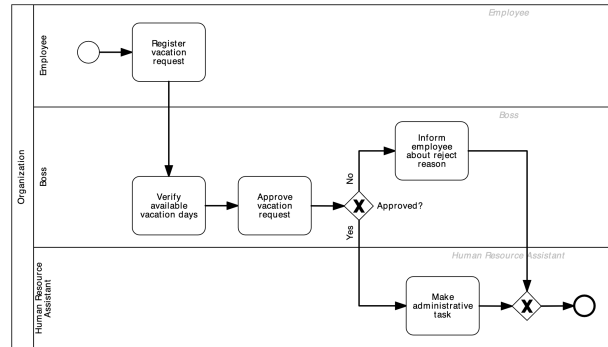


Figure 1: BPMN process specification

angles) performed by different resources or actors (lanes), ii) the sequence flow (directed edges), iii) the gateways that depict parallel, exclusive, and complex forking or joining of paths, and iv) events indicating triggers such as start or end of the process.

A common challenge of using such a specification is the ability to interpret complex specifications, as the process of applying for a vacation will be used by all employees in the organization, who may not have the expertise to interpret the specification. Hence, organizations often provide text-based descriptions that describe the steps of a business process in natural language. Prior work has focused on transforming process models into an intuitive textual description (Leopold et al., 2014), suitable for business process professionals. Yet, there are challenges in navigating and interpreting appropriate textual descriptions through a large repository of process artifacts (van der Aa et al., 2015). One widely plausible approach that organizations have explored to address the challenge of navigating through a large number of documents is to provide chatbots that provide responses based on information present in a knowledge base.

Recent advances have chatbots supported by LLMs thus enabling users ask questions in natural language and receive relevant responses grounded

070 on domain-specific data <sup>1,2</sup>. While engaging in  
071 multi-turn conversations with humans is a fun-  
072 damental capability of LLMs, there has been no  
073 quantitative evaluation of dialogues generated for a  
074 multi-turn conversation based on a process model.  
075 To the best of our knowledge, there is no proper  
076 data set that can be used for this purpose. In this  
077 work, we make the following contributions:

- 078 • Provide a dialog data set grounded on process  
079 descriptions and model specification (Process-  
080 Chat Dataset)
- 081 • Provide a baseline solution for the problem  
082 of end-to-end training of a process model  
083 grounded chatbot and evaluate it in a zero-  
084 shot setting.

## 085 2 Related Work

086 Quantitative evaluation of chatbots for process man-  
087 agement requires datasets depicting process knowl-  
088 edge. There are datasets that are annotated to en-  
089 able the extraction of business artifacts from natural  
090 language descriptions. The most recent is the PET  
091 dataset (Bellan et al., 2022), consisting of a corpus  
092 of business process descriptions annotated with ac-  
093 tivities, gateways, actors, and flow information of  
094 45 processes. Another dataset of 17 process de-  
095 scriptions annotated with declarative constraints  
096 depicting the relation between various activities  
097 has been evaluated and provided by (van der Aa  
098 et al., 2019). A large dataset of 73 process model  
099 descriptions and their BPMN specifications has  
100 been released (Sánchez-Ferreres et al., 2018). The  
101 authors provide ground truth references to missing  
102 alignments between the process description and the  
103 process model. The alignment constraints refer to:  
104 i) cardinality: each activity is aligned to exactly one  
105 sentence, and each gateway is aligned to exactly  
106 one sentence or none at all, ii) ordering: sentences  
107 referring to two activities in the textual description  
108 follow the same order as the process model. A sim-  
109 ilar ordering constraint is also defined for parallel  
110 and exclusive gateways. We use the dataset of these  
111 73 process descriptions to build the dialog dataset  
112 for evaluating chatbots for business processes.

<sup>1</sup><https://www.salesforce.com/news/press-releases/2023/09/12/salesforce-platform-news-dreamforce/>

<sup>2</sup><https://blogs.microsoft.com/blog/2023/03/16/introducing-microsoft-365-copilot-your-copilot-for-work/>

## 3 The ProcessChat Dataset 113

114 In this section, we describe the proposed dataset,  
115 ProcessChat, which includes natural language con-  
116 versations over BPMN processes. To construct  
117 this dataset, we rely on 73 BPMN processes and  
118 natural language process descriptions provided by  
119 (Sánchez-Ferreres et al., 2018) that have been ob-  
120 tained from 11 different industrial and academic  
121 sources<sup>3</sup>. While the original purpose of this repos-  
122 itory was to publish a tool to compute an align-  
123 ment between BPMN process models and natural  
124 language descriptions, we simply use the process  
125 representations in terms of both BPMN and natural  
126 language to build our chat dataset. We start with  
127 constructing multiple dialog flows for each process  
128 using BPMN. While we focus on BPMN, the ap-  
129 proach presented in this paper is independent of the  
130 specific notation used to define a process model.

### 3.1 Dialog grounded on process knowledge 131

132 A chatbot, grounded on process knowledge, con-  
133 siders three different kinds of nodes in BPMN (ex-  
134 ample in Figure 2): i) Events, ii) Activities, and  
135 iii) Gateways. Events represent something happen-  
136 ing in the process. Common events in a process  
137 are the start and end events depicting the start and  
138 end of the process. Additionally, messages and  
139 timers are depicted as intermediate events in a pro-  
140 cess. Activities represent steps in the process that  
141 are performed. Gateways are used to control the  
142 process flows. The two most commonly depicted  
143 gateways are: i) an exclusive gateway (Decision)  
144 which is used to create alternative paths within a  
145 process flow, and an exclusive join gateway which  
146 is used to merge alternative paths, ii) a parallel  
147 gateway having parallel outgoing process flows,  
148 and a parallel join gateway that waits for all in-  
149 coming flows before triggering the next step in the  
150 sequence. Further, the model has pools which are a  
151 graphical representation of a participant or an actor  
152 responsible for the execution of the activities.

153 The characteristics of the process-grounded  
154 dataset are to provide the following kinds of in-  
155 teractions to the user as detailed by (López et al.,  
156 2019):

- 157 • Guide the user step-by-step through the pro-  
158 cess by providing details of the next step to  
159 execute having the context of the previous  
160 steps the user has been guided through.

<sup>3</sup>[https://github.com/setzer22/alignment\\_model\\_text](https://github.com/setzer22/alignment_model_text)

- 161 • Ask the user to provide information when  
162 the process reaches an exclusive gateway to  
163 choose the relevant path and guide the user to  
164 the appropriate path
- 165 • Inform the user about the actor involved in  
166 performing the task
- 167 • Provide the user with the list of activities per-  
168 formed by an actor
- 169 • Inform the user on reaching the end of the  
170 process.

171 An example of a partial conversation flow based  
172 on a Court Hearing process as shown in the BPMN  
173 diagram in Figure 2 is given below.

- 174 T1 **User:** What is the first step of the process for Court  
175 Hearing?
- 176 T1 **Assistant:** The Organization will check the files
- 177 T2 **User:** What is the next step?
- 178 T2 **Assistant:** The Organization wonders if files are miss-  
179 ing
- 180 T3 **User:** What is the next step?
- 181 T3 **Assistant:** What is the outcome of Files missing?
- 182 T4 **User:** Yes
- 183 T4 **Assistant:** The Organization will initiate a search
- 184 T5 **User:** What are the responsibilities of Organization?
- 185 T5 **Assistant:** Check files, Files missing?, Initiate search,  
186 Track physical files, Hand files to associate, Distribute  
187 judge’s lawlist, Conduct directions hearings
- 188 T7 **User:** What is the next step?
- 189 T7 **Assistant:** The Organization will track physical files
- 190 T8 **User:** What is the next step after Track physical files
- 191 T8 **Assistant:** The Organization will hand files to the asso-  
192 ciate and The Organization will distribute the judge’s  
193 lawlist
- 194 T9 **User:** What is the next step after Hand files to associate
- 195 T9 **Assistant:** The Organization will distribute the judge’s  
196 lawlist
- 197 T10 **User:** What is the next step?
- 198 T10 **Assistant:** The Organization will conduct hearings to  
199 give directions
- 200 T11 **User:** What is the next step?
- 201 T11 **Assistant:** You have reached the end of the process

202 The dialog starts by assuming that a user will  
203 initiate the conversation by asking about the first  
204 step in a given process. Then the assistant will  
205 respond by describing the first step. The user is  
206 expected to step through the process by repeatedly  
207 asking “What is the next step?” or could ask a  
208 specific question asking “What is the next step  
209 after A?” (A being an activity in the process), and  
210 occasionally asking role-based questions like “Who  
211 is responsible for this step?” and “What are the  
212 responsibilities of so and so?”. When the process  
213 is finished, the assistant will inform the user that  
214 the end of the process is reached. Each user query  
215 and an assistant response is a ‘turn’.

216 In the case of an exclusive gateway, multiple  
217 next steps are possible following a given step but

218 are dependent on some conditions. The assistant  
219 will then pose a query to the user, asking for the  
220 response to the question specified in the exclusive  
221 gateway (at T3 in the conversation). Depending on  
222 the user’s response, the assistant will specify the  
223 next step.

224 In the case of a parallel gateway, multiple next  
225 steps are possible following a given step, and the  
226 assistant will respond by listing all of them (at T8  
227 in the conversation). It will then step through all  
228 the parallel paths based on the user choosing a path  
229 and enquiring the next step (at T9 enquiring about  
230 the next step after *hand files to associate*). There  
231 is another dialog flow created where the user can  
232 enquire on the other parallel path (“What is the  
233 next step after *distribute judge’s lawsuit?*”)

234 Timer or other intermediate events are presented  
235 to the user as the next steps. For example, if there  
236 is a timer event “design complete”, the response of  
237 the chatbot to “What is the next step?” will be “The  
238 next step is to get the design completed.” Hence,  
239 the process chat considers events, activities, gate-  
240 ways, and participants (or resources) to construct  
241 the dialog.

### 242 3.2 Dialog Data construction

243 We systematically iterate over the paths in the  
244 BPMN process model and for each path, construct  
245 a dialog flow. Each dialog flow consists of three  
246 parts: i) process description in natural language  
247 (NL), ii) process description in constrained natural  
248 language (CNL), and iii) process path traversal for  
249 the multi-turn dialog flow.

250 **Process Description in NL:** We envision a scenario  
251 where an LLM-driven chatbot (or *assistant*)  
252 will be able to answer questions on the process  
253 given a textual representation of the process. The  
254 textual descriptions in the dataset have been pro-  
255 vided by experts covering different styles a variety  
256 of styles, The creation of descriptions by experts  
257 included three steps (i) Study the process model di-  
258 agram. (ii) Write the textual description. (iii) Com-  
259 pare the textual description with the process model  
260 to and make sure the text accurately describes the  
261 process model. This final step aimed to reduce  
262 the amount of inconsistencies between texts and  
263 models.

264 For example, the provided NL description for  
265 the example process is as follows: “*Each morning,*  
266 *the files which have yet to be processed need to*  
267 *be checked, to make sure they are in order for the*

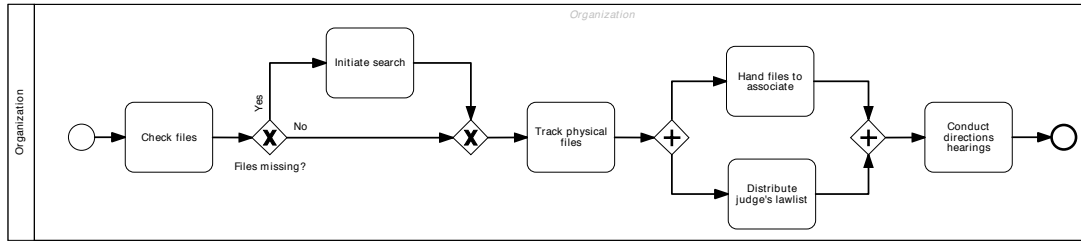


Figure 2: Court Hearing process

268 *court hearing that day. If some files are missing,*  
 269 *a search is initiated, otherwise the files can be*  
 270 *physically tracked to the intended location. Once*  
 271 *all the files are ready, these are handed to the*  
 272 *Associate, and meantime the Judges Lawlist is*  
 273 *distributed to the relevant people. Afterward, the*  
 274 *directions hearings are conducted.”*  
 275

276 **Process Description in CNL:** BPMN models  
 277 can provide precise information as compared to  
 278 textual descriptions. We would also want to know  
 279 if an LLM could reason and answer questions with  
 280 a BPMN model as an input. However, as discussed  
 281 in prior work (Grohs et al., 2023), as BPMN specifi-  
 282 cation uses the XML format, it can be verbose. The  
 283 example BPMN of Figure 2, having 6 activities, 2  
 284 events and 4 gateways results in 2040 tokens of  
 285 an open-source LLM such as Mistral (Jiang et al.,  
 286 2023). The NL description for the equivalent spec-  
 287 ification results in 94 tokens. Hence, similar to the  
 288 approach proposed by (Grohs et al., 2023), we syn-  
 289 thesize a CNL by parsing the BPMN to represent  
 290 the process flow. Using a constrained language  
 291 to represent the flow reduces ambiguity. CNL is  
 292 NL with restricted grammar and hence simplified  
 293 and standardized sentence structure. We hypothe-  
 294 size that CNL would be precise as it represents the  
 295 BPMN model and yet results in improved semantic  
 296 quality for an LLM to reason. There are constraints  
 297 the CNL represents: i) an ordering constraint such  
 298 as *step 1* must happen before *step 2* where *step 1*  
 299 and *step 2* are consecutive steps in a path, and ii)  
 300 conditional constraints for an XOR gateway such as  
 301 if the response to *query 1* is Yes, then execute *step*  
 302 *3*, where *step 3* is a process step which immedi-  
 303 ately follows an exclusive gateway in the BPMN  
 304 represented by *query 1*.

305 The auto-generated CNL rules for the BPMN  
 306 model (Figure 2) are:

- 307 • if the response to Files missing? is Yes then  
 308 Initiate search

- Distribute judge’s lawlist must happen before  
 Conduct directions hearings
- Initiate search must happen before Track phys-  
 ical files
- Track physical files must happen before Hand  
 files to associate
- Hand files to associate must happen before  
 Distribute judge’s lawlist
- Check files must happen before Initiate search

In case of a loopback from a process step B to  
 a previous process step A, while A must always  
 happen before B in any execution path through the  
 process, B only happens before A at the end of the  
 loop. So in such cases, we add the CNL rule “B  
 can happen before A”. For the same court hearing  
 process, use of CNL reduces the number of tokens  
 to 142.

**Process Path Traversal for Dialog flow:** The  
 dialog flow  $d$  between a user  $u$  and the chatbot  
 (or assistant)  $a$ , is represented as a sequence of  
 utterances  $d = \{c_1^u, c_1^a, c_2^u, c_2^a \dots c_m^u, c_m^a\}$ , where  
 $m$  denotes the number of exchanges or turns in  
 the dialog. The BPMN can be considered as a set  
 of element nodes  $N$  and edges  $E$ . As a path in  
 the BPMN is traversed, for each element of type  
 activity, a user question and an assistant response  
 describing the activity are created. For each ex-  
 clusive gateway, an agent question is associated  
 with the gateway, and the response from the user  
 for the chosen path is created. Hence, an exclu-  
 sive gateway always results in an assistant asking  
 a question. In case the element type is a parallel  
 gateway, we first create a user query and an agent  
 response that reflects the first activities on each  
 of the parallel paths. We then navigate each valid  
 path with the user specifically choosing the first  
 tasks in each of the parallel branches as a query.  
 We thus create multiple dialog flows for each pro-  
 cess having exclusive and parallel gateways. Once  
 all the steps are complete, we create a user query  
 and assistant response depicting the end of the pro-



cess. The objective of using the data is to learn the next response, which takes (i) the dialog history  $h = \{c_1^u, c_1^a \dots, c_u^i\}$ , and (ii) a process description (as NL or CNL) and predict the agent response ( $y = c_i^a$ ).

Further, we generate the assistant response for each element in BPMN using an LLM (Mistral (Jiang et al., 2023)). The activity label *Conduct direction hearings*, and the participant name (*Organization*) is used to generate the utterance: ‘*The Organization will conduct hearings to give directions*’. A title for each process is also generated using an LLM by providing the process description as input. For example, when the user starts a conversation, the user utterance we generate starts as *what is the first step in the process for <process title>?* In the future, it is possible to use the ProcessChat data to baseline a retrieval augmented generation framework (Gao et al., 2024) by retrieving relevant process descriptions to generate assistant responses.

### 3.3 ProcessChat Data

The ProcessChat<sup>4</sup> data is divided into train, test, and validation splits at the process level so that processes are not shared across splits. Each process generates multiple conversation flows depending on the number of exclusive gateways and the paths taken by the user. Again, each conversation flow generates multiple samples where each sample denotes the expected response of the assistant in a single turn given the conversation context. The statistics of the dataset, including processes, conversation flows, samples, and various question types are given in Table 1. Next-step questions include user queries related to the steps of the process (first or next), whereas resource questions consist of role and responsibility-related queries. Assistant questions are generated to identify the outcome of decision points (exclusive gateways), which in turn will be asked by the assistant to decide the next step.

## 4 Baseline System

We aim to investigate the performance of LLMs on process dialog generation tasks, with a particular focus on decoder-based LLMs. In our baseline experiments, we explore zero-shot prompting, few-shot learning, and fine-tuning. For the fine-tuning

Data Split	train	val	test
Processes	51	7	15
Conversation Flows	230	24	62
Samples	3092	439	846
Next-step Questions	2081	290	551
Resource Questions	665	107	182
Assistant Questions	355	42	122

Table 1: Dataset statistics

process, we employ the LoRA method (a Low-Rank Adapter framework) (Hu et al., 2021). To the best of our knowledge, this is the first instance of LoRA being applied to fine-tuning LLMs for process-grounded dialog tasks.

### 4.1 Prompting

We investigate a varied collection of pre-trained LLMs, all based on the transformer architecture (Vaswani et al., 2017). This collection comprises three distinct LLMs, each trained on instruction data and available in multiple versions with varying parameter sizes. This results in models with parameter sizes ranging from 7 billion to 47 billion.

**Mistral-7B-Instruct-v0.2 (Jiang et al., 2023):** The model is a 7.3B parameter which is an instruction-tuned version of the base Mistral 7B model. It has outperformed larger models in multiple NLP benchmarks.

**granite-13b-chat-v2 (IBM, 2024):** The model was initialized from a base model trained on 1.25 trillion tokens and also relies on synthetic data that is designed to improve the model’s conversational, safety, and instruction following capabilities. The model is aligned to chat instructions including a conversation history.

**Mixtral-8x7B-Instruct-v0.1 (Jiang et al., 2024):** This model from Mistral AI is a “mixture of experts” model which is a decoder-only model where the feedforward block picks from a set of 8 distinct groups of parameters allowing it to use only a fraction of the total set of parameters per token. It has 46.7B total parameters but only uses 12.9B parameters per token.

These models were chosen primarily because they support a large (around 8K) input token size. To compare the results from these models with fine-tuning based results we considered the test set only by leaving the training and validation splits. We use few-shot learning, also denoted as K-shot, with

<sup>4</sup><https://anonymous.4open.science/r/ProcessChat-4FF4/README.md>

K representing the number of examples provided, where in our case, examples are randomly sampled from the training set. To understand the effect of the number of few-shot examples we experimented with 0, 1, 3, and 5 examples in the prompt. The prompt we used is similar to what is shown in Figure 4 included in Appendix A, where some few-shot examples consisting of conversations represented in the same format are added before the current conversation.

## 4.2 Fine-tuning

We experimented with fine-tuning two LLMs, IBM’s granite-13b-base-v1 (IBM, 2024) and the open-source Mistral-7B-v0.1 (Jiang et al., 2023) from Mistral AI<sup>5</sup>. To conserve GPU memory, fine-tuning is done in a parameter-efficient manner using LoRA (Hu et al., 2021) with a rank of 4. LoRA (short for Low-Rank Adaptation) is a technique that freezes the pre-trained model weights and injects trainable rank decomposition matrices into each layer of the Transformer (Vaswani et al., 2017) architecture, thereby greatly reducing the number of trainable parameters for fine-tuning. We chose this method because fully fine-tuning such LLMs is expensive and can lead to loss of generalizability due to catastrophic forgetting. The models are trained on the samples in the training set and tested on the samples in the test set, while the samples in the validation set are used for evaluating model performance during training. The models are trained for 3 epochs with a batch size of 2 on a single NVIDIA A100 GPU with 80 GB memory. We use the latest checkpoint as the final one for testing the models. We spent around 12 GPU hours for fine-tuning granite-13b-base-v1 and around 8 GPU hours for fine-tuning Mistral-7B-v0.1 (so the total computational budget was around 20 GPU hours).

For each turn of the assistant in the dataset, we obtain a sample data point where the input to the model is a prompt consisting of a system instruction along with the process description and the dialogue history, and the output is the response given by the assistant to the user. For the Court Hearing process described above, a set of example input prompts and output responses are shown in Appendix A.

<sup>5</sup><https://mistral.ai>

## 4.3 Evaluation and Baseline Results

We considered three settings based on the inclusion of process description for evaluating the ability of the models (as assistants) to generate responses:

1. Including Both NL & CNL descriptions
2. Including only NL description
3. Including only CNL description

We compare the assistant response generated by the model with the expected response of the assistant at each turn in the test set using ROUGE-L (Lin, 2004) and BLEU (Papineni et al., 2002) scores. We then compute average ROUGE-L and BLEU scores for all the samples in the test set to get the overall performance of the model. We also separate out samples where the user asks resource (i.e. role/responsibility) based questions from the test set and measure the average ROUGE-L and BLEU scores for only those samples. This gives us an idea of how good the model is at identifying the roles of each resource in the process. Furthermore, we check whether the model is generating questions at the right turn in the dialogue or not. To do this, we compute two additional metrics on the test set: (i) Question Precision (Q-Pr) which measures how often a generated question is expected at that turn, and (ii) Question Recall (Q-Re) which measures how often a question is generated when it is expected at that turn.

### 4.3.1 Prompt-based results

From Table 2, it is evident that more few-shot examples always help the model to achieve better scores. We see that the Mixtral-8x7B-Instruct-v0.1 model is outperforming the other two models in all metrics. Also, from the results we see that using CNL descriptions, with or without NL descriptions, has a definite advantage over using only NL descriptions in terms of ROUGE-L and BLEU scores. However, it is unclear whether using CNL descriptions in conjunction with NL descriptions is better than only using CNL descriptions or vice versa. We also notice that these models are better at answering resource-related questions in general. While Mixtral-8x7B-Instruct-v0.1 can show good Question Precision, all the prompt-based models show poor Question Recall. This may be because the models are not able to understand the process well enough to know when a question must be asked.

Interestingly, the performance of Mistral-7B-Instruct-v0.2 and Mixtral-8x7B-Instruct-v0.1 in the five-shot setting on the resource-related questions

Expt.	Model	Policy	ROUGE-L (overall)	BLEU (over- all)	ROUGE-L (resource)	BLEU (re- source)	Q-Pr	Q-Re
LoRA	granite-13b-base-v1	NL+CNL	<b>0.803</b>	<b>0.560</b>	<b>0.808</b>	<b>0.312</b>	0.745	0.860
	granite-13b-base-v1	NL	0.654	0.360	0.779	0.244	0.617	0.647
	granite-13b-base-v1	CNL	0.758	0.494	0.799	0.294	<b>0.781</b>	<b>0.877</b>
	Mistral-7B-v0.1	NL+CNL	0.497	0.221	0.492	0.116	0.172	0.180
	Mistral-7B-v0.1	NL	0.466	0.181	0.495	0.122	0.135	0.106
	Mistral-7B-v0.1	CNL	0.482	0.215	0.492	0.120	0.226	0.238
5-shot Prompt	granite-13b-chat-v2	NL+CNL	0.522	0.237	0.602	0.207	0.200	0.008
	granite-13b-chat-v2	NL	0.457	0.169	0.570	0.157	0.050	0.008
	granite-13b-chat-v2	CNL	0.452	0.207	0.396	0.139	0.133	0.016
	Mistral-7B-Instruct-v0.2	NL+CNL	0.562	0.306	0.544	0.177	0.000	0.000
	Mistral-7B-Instruct-v0.2	NL	0.511	0.222	0.567	0.154	0.333	0.008
	Mistral-7B-Instruct-v0.2	CNL	0.577	0.307	0.645	0.204	0.250	0.008
	Mixtral-8x7B-Instruct-v0.1	NL+CNL	0.653	<b>0.379</b>	0.757	0.284	0.823	0.114
	Mixtral-8x7B-Instruct-v0.1	NL	0.577	0.274	0.732	0.218	<b>1.000</b>	0.033
	Mixtral-8x7B-Instruct-v0.1	CNL	<b>0.656</b>	0.363	<b>0.809</b>	<b>0.286</b>	0.667	<b>0.147</b>
3-shot Prompt	granite-13b-chat-v2	NL+CNL	0.509	0.217	0.592	0.215	0.222	0.049
	granite-13b-chat-v2	NL	0.456	0.163	0.567	0.152	0.204	0.082
	granite-13b-chat-v2	CNL	0.463	0.187	0.554	0.195	0.111	0.016
	Mistral-7B-Instruct-v0.2	NL+CNL	0.549	0.301	0.482	0.160	0.636	0.057
	Mistral-7B-Instruct-v0.2	NL	0.522	0.211	0.642	0.173	1.000	0.008
	Mistral-7B-Instruct-v0.2	CNL	0.578	0.290	0.705	0.226	0.333	0.024
	Mixtral-8x7B-Instruct-v0.1	NL+CNL	0.651	0.355	0.785	0.289	0.692	0.147
	Mixtral-8x7B-Instruct-v0.1	NL	0.573	0.252	0.761	0.223	0.857	0.049
	Mixtral-8x7B-Instruct-v0.1	CNL	0.634	0.335	0.797	0.280	0.695	0.131
1-shot Prompt	granite-13b-chat-v2	NL+CNL	0.458	0.183	0.124	0.497	0.000	0.000
	granite-13b-chat-v2	NL	0.386	0.468	0.113	0.131	0.333	0.016
	granite-13b-chat-v2	CNL	0.457	0.188	0.531	0.151	0.000	0.000
	Mistral-7B-Instruct-v0.2	NL+CNL	0.489	0.238	0.374	0.087	0.000	0.000
	Mistral-7B-Instruct-v0.2	NL	0.478	0.194	0.552	0.142	0.000	0.000
	Mistral-7B-Instruct-v0.2	CNL	0.479	0.231	0.407	0.091	0.000	0.000
	Mixtral-8x7B-Instruct-v0.1	NL+CNL	0.577	0.283	0.600	0.169	0.750	0.024
	Mixtral-8x7B-Instruct-v0.1	NL	0.496	0.205	0.563	0.148	0.000	0.000
	Mixtral-8x7B-Instruct-v0.1	CNL	0.581	0.287	0.620	0.157	0.500	0.016
0-shot Prompt	granite-13b-chat-v2	NL+CNL	0.057	0.012	0.015	0.055	0.000	0.000
	granite-13b-chat-v2	NL	0.071	0.063	0.019	0.017	0.000	0.000
	granite-13b-chat-v2	CNL	0.044	0.010	0.052	0.021	0.000	0.000
	Mistral-7B-Instruct-v0.2	NL+CNL	0.388	0.454	0.102	0.163	0.000	0.000
	Mistral-7B-Instruct-v0.2	NL	0.371	0.123	0.428	0.094	0.000	0.000
	Mistral-7B-Instruct-v0.2	CNL	0.390	0.175	0.458	0.104	0.200	0.008
	Mixtral-8x7B-Instruct-v0.1	NL+CNL	0.192	0.113	0.175	0.052	0.000	0.000
	Mixtral-8x7B-Instruct-v0.1	NL	0.268	0.109	0.236	0.054	0.000	0.000
	Mixtral-8x7B-Instruct-v0.1	CNL	0.068	0.036	0.118	0.043	0.000	0.000

Table 2: Results of various baseline approaches on the test set. The best scores obtained by the fine-tuned and prompt-based models are marked in bold.

using only CNL descriptions is better than the performance using only NL or a combination of NL and CNL descriptions. We surmise that this may be because the model gets a compact description of the process in CNL which helps it to infer the answers to the resource-related questions in a succinct way using the dialog context.

### 4.3.2 Fine-tuning results

From Table 2 we see that the fine-tuned granite-13b-base-v1 is better than the fine-tuned Mistral-7B-v0.1 by a large margin. In fact, the fine-tuned Mistral-7B-v0.1 seems to be performing worse than the prompt-based Mistral-7B-Instruct-v0.2 model. This may be because of higher quality instruction tuning done by Mistral AI which makes the model very robust to various types of human instructions<sup>6</sup>. However, in terms of ROUGE-L and BLEU scores, the fine-tuned granite-13b-base-v1 is better than all the prompt-based models. We see that using both NL and CNL descriptions is better than only using CNL descriptions, which in turn is better than only using NL descriptions. This may be because fine-tuning allows the model to get useful information from both NL and CNL descriptions which may be complementary to each other. While the fine-tuned granite-13b-base-v1 may have lower Question Precision than some of the prompt-based models, it has the highest Question Recall among all the models. Interestingly, using CNL descriptions alone achieves higher Question Precision and Question Recall than using only NL or both NL and CNL descriptions. This may be because the CNL description indicates the decision points, which may not be indicated in the NL description.

As observed during prompting, the performance of both the models on the resource-related questions using only CNL descriptions is not much worse than the performance using only NL or a combination of NL and CNL descriptions. In the case of granite-13b-base-v1, it is actually better than using only NL descriptions.

Figure 3 shows a graphical representation of the scores attained by the models (considering only fine-tuning and five-shot prompting with both NL and CNL descriptions), to compare average ROUGE-L scores for the samples in the test set for the overall assistant responses to all types of user questions, the responses related to questions on the

<sup>6</sup>Notably, fine-tuning Mistral-7B-Instruct-v0.2 itself using LoRA yields even poorer results, presumably because the effect of instruction tuning is destroyed.

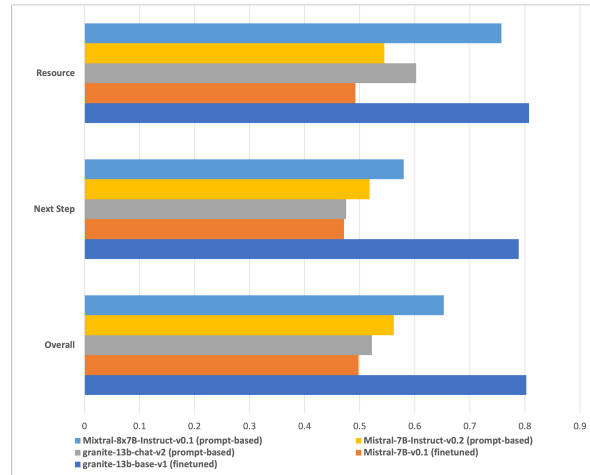


Figure 3: ROUGE-L scores of answers provided by the system for various kinds of user queries (overall, resource-related, and next-step related), considering only fine-tuning and five-shot prompting for the setting where both NL and CNL process descriptions are used.

next step, and the responses related to questions on a resource. The fine-tuned model performed better in all three scenarios, while the prompt-based models did well for only resource-related queries.

## 5 Conclusion and Future Work

In this work, we presented a view of providing chatbot interactions with grounding on process knowledge. We conduct experimental results to verify the suitability of LLMs to infer and provide guidance based on process descriptions in natural language and constrained natural language. We generate and release the ProcessChat dataset, which contains 316 dialog flows grounded on 73 process models. We propose the baseline solutions that evaluate a prompt-based model and a parameter-efficient fine-tuned model using LoRA. Our baseline results show significant improvement using a fine-tuned model with a small training dataset of 230 dialog flows. We foresee multiple directions for future research: i) The ProcessChat dataset is suited for step-by-step guidance. It can be extended to provide an interaction where the user should can ask any specific query regarding the process. This would require generating dialog interactions on any step in the process. ii) The dataset does not consider data flow as the specification had limited information on data artifacts. Incorporating data perspective would ensure a comprehensive experience. iii) Finally, this dataset can be extended for interactions across multiple processes in a repository.



## 6 Limitations

In the process notation, splitting parallel gateway multiplies the incoming sequence flow into several outgoing sequence flows that run simultaneously. A joining parallel gateway waits for all incoming sequences to terminate before combining them all in one outgoing flow. This leads to several variations in the flow execution. For example, activities in one parallel path can be interspersed with activities in another parallel path. However, To reduce redundancy, we considered the sequence of activities in one flow intact. Hence, in our dataset, a process with  $n$  parallel flows will lead to only  $n$  output flows, where a user completes dialog interaction of one parallel path and then proceeds to the next path.

Another limitation in the generation of the dataset is the use of an LLM to form a assistant response or question from the BPMN specification. As the specification only contains activities names and the actors performing it, we generate the assistant response using **Mistral-7B-Instruct-v0.2 (Jiang et al., 2023)**. In the conversation presented for the court hearing process, assistant responses at T1, T2 or the question at T3 was generated by the LLM. Hence, the assistant response or the ground truth response is dependent on the LLM we have used. However, we have manually evaluated these generated responses.

## References

Patrizio Bellan, Han van der Aa, Mauro Dragoni, Chiara Ghidini, and Simone Paolo Ponzetto. 2022. PET: an annotated dataset for process extraction from natural language text tasks. In *BPM 2022 International Workshops, Münster*, volume 460 of *LNBIP*, pages 315–321.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. 2024. [Retrieval-augmented generation for large language models: A survey](#). *Preprint*, arXiv:2312.10997.

Michael Grohs, Luka Abb, Nourhan Elsayed, and Jana-Rebecca Rehse. 2023. Large language models can accomplish business process management tasks. In *BPM 2023 International Workshops, Utrecht*, volume 492 of *LNBIP*, pages 453–465.

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

IBM. 2024. Granite foundation models. <https://www.ibm.com/downloads/cas/X9W406BM>.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

Henrik Leopold, Jan Mendling, and Artem Polyvyanyy. 2014. Supporting process model validation through natural language generation. *IEEE Trans. Software Eng.*, 40(8):818–840.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Anselmo López, Josep Sànchez-Ferreres, Josep Carmona, and Lluís Padró. 2019. From process models to chatbots. In *CAiSE 2019, Rome, Proceedings*, volume 11483 of *LNCIS*, pages 383–398.

OMG. 2011. Business Process Model and Notation (BPMN), Version 2.0. *Object Management Group*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Josep Sànchez-Ferreres, Han van der Aa, Josep Carmona, and Lluís Padró. 2018. Aligning textual and model-based process descriptions. *Data & Knowledge Engineering*, 118:25–40.

Han van der Aa, Claudio Di Ciccio, Henrik Leopold, and Hajo A. Reijers. 2019. Extracting declarative process models from natural language. In *CAiSE 2019, Rome, 2019, Proceedings*, volume 11483 of *LNCIS*, pages 365–382.

Han van der Aa, Henrik Leopold, and Hajo A. Reijers. 2015. Detecting inconsistencies between process models and textual descriptions. In *BPM - 13th International Conference, Austria, 2015, Proceedings*, volume 9253, pages 90–105.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

## A Appendix

Set of example input prompts and output responses are given by Figure 4, where each input terminates with the user’s utterance and the expected output is given by the assistant’s response.

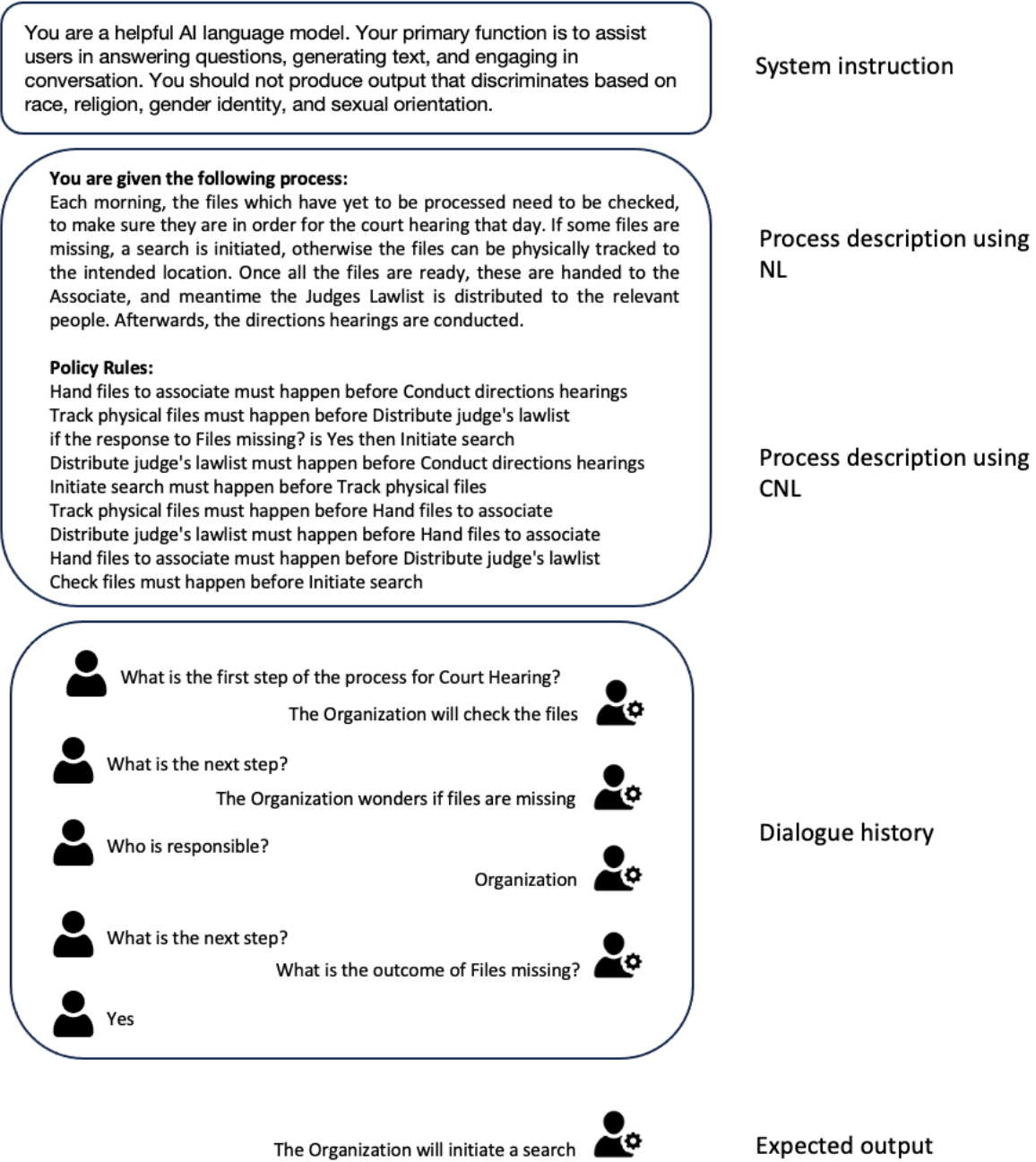


Figure 4: Example input prompt and expected response