

PAGED: A Benchmark for Procedural Graphs Extraction from Documents

Weihong Du¹² Wenrui Liao¹² Hongru Liang^{12*} Wenqiang Lei¹²

¹College of Computer Science, Sichuan University, China

²Engineering Research Center of Machine Learning and Industry Intelligence,
Ministry of Education, China

{duweihong, liaowenrui}@stu.scu.edu.cn

{lianghongru, wenqianglei}@scu.edu.cn

Abstract

Automatic extraction of procedural graphs from documents creates a low-cost way for users to easily understand a complex procedure by skimming visual graphs. Despite the progress in recent studies, it remains unanswered: *whether the existing studies have well solved this task (Q1)* and *whether the emerging large language models (LLMs) can bring new opportunities to this task (Q2)*. To this end, we propose a new benchmark PAGED, equipped with a large high-quality dataset and standard evaluations. It investigates five state-of-the-art baselines, revealing that they fail to extract optimal procedural graphs well because of their heavy reliance on hand-written rules and limited available data. We further involve three advanced LLMs in PAGED and enhance them with a novel self-refine strategy. The results point out the advantages of LLMs in identifying textual elements and their gaps in building logical structures. We hope PAGED can serve as a major landmark for automatic procedural graph extraction and the investigations in PAGED can provide valuable insights into the research on logical reasoning among non-sequential elements. The code and dataset are available in <https://github.com/SCUNLP/PAGED>.

1 Introduction

Procedural graphs, though can intuitively represent the execution of actions for goal achievement (Mouchi, 1980; Ren et al., 2023), suffer from the high cost of expert-construction (Herbst and Karagiannis, 1999; Maqbool et al., 2019). The automatic extraction of procedural graphs from procedural documents thus has huge potential, as it would enable users to easily understand how to logically perform a goal (e.g., how a restaurant serves the customers) by skimming visual graphs (e.g., Figure 1(a) instead of reading lengthy documents (e.g., Figure 1(b)).

However, obtaining optimal procedural graphs is not easy — as shown in Figure 1(a), it requires representing not only sequential actions in the procedure (e.g., I-4 → I-5), but also non-sequentially executed actions (e.g., I-7.1 & I-7.2) and vital constraints for the actions (e.g., C-2). Off-the-shelf attempts only meet part of the requirements. For example, Bellan et al. (2023); Ren et al. (2023) fail to represent “customers can order only the dishes or drinks, or both” (I-2.1 & I-2.2) due to its inherent limitation for representing complex non-sequential actions. Besides, current solutions only focus on hand-written rules (Sholih et al., 2022) or customized networks (Bellan et al., 2023) on a small group of cherry-picked instances. This brings up the question of *whether the existing studies have well solved the automated extraction of procedural graphs from procedural documents (Q1)*. If the answer is “no”, we are also interested in *whether the emerging large language models (LLMs) can bring new opportunities to this task (Q2)*.

To answer Q1, we propose to construct a standard benchmark for the Procedural Graphs Extraction from Documents (PAGED). As far as we know, there lack of large-scale datasets of document-graph pairs for training and evaluating optimal procedural graph extraction models (cf., Table 1). We want to equip PAGED with the largest publicly available dataset. Although there exist plenty of procedural documents on the Internet, it is too costly to filter the low-quality ones and annotate optimal procedural graphs matching all requirements. As a remedy, we build the dataset based on a model collection of business process (Dumas et al., 2018), which has summarized business processes into high-quality procedural graphs with complete sequential actions, non-sequential actions, and constraints. Thus, constructing procedural document-graph pairs turns into assigning a suitable procedural document to a given procedural graph. We approach it via a three-stage pipeline — we pro-

* Corresponding author

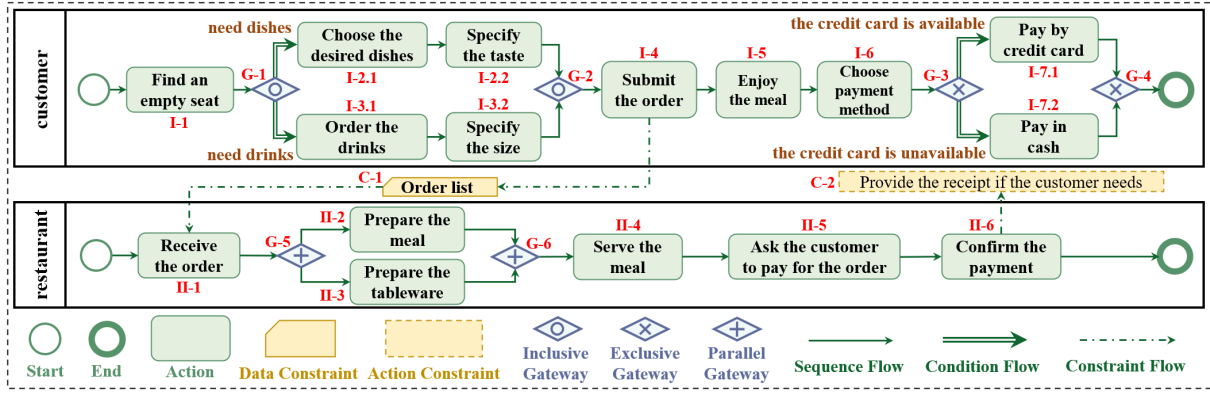


Figure 1: The procedure of how a restaurant serves the customers in procedural graph (a) and document (b).

Firstly, the customer needs to find an empty seat. If the customer needs dishes, then choose the desired dishes and specify the taste. If the customer needs drinks, then order the drinks and specify the size. The customer then submits the order, which is added to the order list. After enjoying the meal, the customer should choose payment method. If the credit card is available, the customer pays by credit card; else if the credit card is unavailable, the customer should pay in cash. For the restaurant, once receiving the order from the order list, it prepares the meal according to the order and prepares the tableware for the customer at the same time. The meal is then served for the customer to enjoy. After that, the restaurant asks the customer to pay for the order and then confirms the payment. Note that the restaurant should provide the receipt if the customer needs.

Figure 1: The procedure of how a restaurant serves the customers in procedural graph (a) and document (b).

gressively transfer the structured information on the procedural graph into natural language text, adjust its narration, and improve the coherence and naturalness, finally generating a suitable document. In this way, we develop a dataset with 3,394 high-quality procedural document-graph pairs that are about ten times larger than the previous largest datasets (Ackermann et al., 2021; Qian et al., 2020). According to the underlying structure of optimal procedural graphs, we introduce three metrics to evaluate five state-of-the-art methods (Sonbol et al., 2023; Neuberger et al., 2023; Sholiq et al., 2022).

To further answer Q2, we investigate the performance of three advanced LLMs (Flan-T5 (Chung et al., 2022), ChatGPT (Ouyang et al., 2022) and Llama2 (Touvron et al., 2023)) and utilize a self-refine strategy to improve the ability of LLMs. In total, we evaluate ten methods in our PAGED benchmark. Extensive experiments on our benchmark reveal that existing studies struggle to accurately extract sequential actions, constraints, and organize non-sequential actions of procedural documents. While LLMs have shown significant improvement in sequential action and constraint extraction, they still face challenges with non-sequential action organization. Our detailed analysis of the results leads us to propose improvement strategies to help large language models better understand non-sequential actions and use correct gateways to represent them. We hope PAGED can be a key mile-

stone for automatic procedural graphs extraction, offering insights into research on logical reasoning among non-sequential elements. In summary, we highlight PAGED as follows:

- We build a novel benchmark named PAGED, which standardly evaluates the progress of current procedural graphs extraction from documents and explores the potential of emerging LLMs.
- We equip PAGED with the largest procedural document-graph dataset, whose high quality is achieved by a three-stage pipeline and verified via both automatic and human evaluation.
- We systematically evaluate state-of-the-art solutions in PAGED and reveal that they have trouble extracting optimal procedural graphs due to their heavy reliance on hand-written rules and limited available data.
- We investigate advanced LLMs in PAGED and empower them with a self-refine strategy, showing their advantages in identifying sequential actions and constraints, and pointing out their gaps in building complex logic of graphs.

2 Related Work

Procedural Graph Extraction Existing studies only meet part of the requirements for optimal procedural graph extraction. Earlier studies mainly

Table 1: Comparisons between our dataset with the existing datasets.

Dataset	Samples Num	Sequential Actions	Non-sequential Actions			Constraints		Publicly Available
			Exclusive Gateway	Inclusive Gateway	Parallel Gateway	Data Constraint	Action Constraint	
Friedrich et al. (2011)	47	✓	✓	✗	✓	✓	✗	✓
Epure et al. (2015)	34	✓	✓	✗	✓	✗	✗	✗
Ferreira et al. (2017)	56	✓	✓	✗	✓	✗	✗	✗
Mendling et al. (2019)	103	✓	✗	✗	✗	✗	✗	✓
Quishpi et al. (2020)	121	✗	✗	✗	✗	✗	✗	Partial
Qian et al. (2020)	360	✗	✗	✗	✗	✗	✗	✓
Ackermann et al. (2021)	358	✗	✗	✗	✗	✗	✗	✓
López et al. (2021)	37	✓	✗	✗	✗	✗	✗	Partial
Bellan et al. (2023)	45	✓	✓	✗	✓	✗	✗	✓
Liang et al. (2023)	200	✓	✗	✗	✗	✗	✗	✓
Ren et al. (2023)	283	✓	✗	✗	✗	✗	✗	✓
ours	3,394	✓	✓	✓	✓	✓	✓	✓

focus on extracting sequential actions (Pal et al., 2021; López et al., 2021; Ren et al., 2023). While Epure et al. (2015); Honkisz et al. (2018); Bellan et al. (2023) explore non-sequential actions, they do not cover scenarios like I-2.1 & I-2.2 in Figure 1(a). Friedrich et al. (2011) aids in data constraint extraction but overlooks action constraints. Besides, current studies heavily rely on hand-written rules and templates (Epure et al., 2015; Honkisz et al., 2018), resulting in poor generalization. Bellan et al. (2023) trains a neural network model but only learns 45 samples, whose effectiveness remains questionable. Hence, we propose to construct a standard benchmark to reveal the performance of existing studies and highlight the challenges for the extraction of optimal procedural graphs.

Datasets As shown in Table 1, current datasets consist of a small group of cherry-picked instances. Some datasets are not publicly available (Epure et al., 2015; Ferreira et al., 2017). Some datasets only focus on sentence-level actions extraction, lacking both sequential and non-sequential actions (Quishpi et al., 2020; Qian et al., 2020; Ackermann et al., 2021). While other studies contain sequential actions Friedrich et al. (2011); Mendling et al. (2019); López et al. (2021); Bellan et al. (2023); Liang et al. (2023); Ren et al. (2023), they do not cover all types of non-sequential actions. Moreover, almost all existing datasets ignore vital constraints related to the actions in procedural graphs. To this end, we construct a new procedural document-graph dataset that is nearly ten times larger than the previous largest datasets (Ackermann et al., 2021; Qian et al., 2020). Each sample consists of a high-quality procedural document and its procedural graph with complete sequential actions, non-sequential actions, and constraints.

Data2Text Data2Text task aims at transferring structured data such as graphs into natural language text (Duong et al., 2023; Lin et al., 2023). Current studies (Su et al., 2021; Kasner and Dušek, 2022) only focus on the transformation of factual knowledge — knowledge about features of things, making it difficult to deal with procedural knowledge — execution of sequential and non-sequential actions in the procedural graphs. Moreover, current studies can only manage discrete components (Ye et al., 2019; Fu et al., 2020), while extracting procedural graphs requires handling complex logic of sequential action, non-sequential actions and their constraints. In this paper, we propose a three-stage pipeline to bridge the gap between complex graphs and lengthy documents, ensure logical descriptions of generated documents, and solve the issues of fluency and coherence in generated documents.

Large Language Models The emerging LLMs have presented competitive results in a wide range of tasks (Zhao et al., 2023), but are barely used for procedural graph extraction. The only exception is Bellan et al. (2022), which makes a shallow attempt to extract sequential actions and deal with partial non-sequential actions with LLMs, and performs poorly for gateway extraction. It remains unanswered whether LLMs’ ability to understand the inherent structure of long contexts can improve the procedural graphs extraction from documents. To this end, we involve Flan-T5 (Chung et al., 2022), ChatGPT (Ouyang et al., 2022) and Llama2 (Touvron et al., 2023) in our PAGED and design a self-refine strategy to demonstrate the opportunities and gaps of LLMs in this task. We hope this can help to explore more possibilities that LLMs bring to this field.

3 Dataset

It is too costly to conduct an expert annotation of optimal procedural graphs for a large number of documents. To this end, we build our dataset upon a model collection of business process (Dumas et al., 2018), which has defined optimal procedural graphs covering the whole business process management lifecycle. The dataset construction turns into a data2text task — generates suitable documents for given procedural graphs.

3.1 Preliminary

Figure 1(a) presents an example of the optimal procedural graph. It describes the procedure of how a restaurant serves customers and involves two actors (customer and restaurant). Each actor starts from the “Start” node, carries out actions following the logic in the graph, and ends at the “End” node. If the actions are performed sequentially, they are connected by the “Sequence Flow”. Otherwise, there is an “Inclusive Gateway”, “Exclusive Gateway” or “Parallel Gateway” indicating the following actions are non-sequential ones. Both the “Inclusive Gateway” (G-1) and “Exclusive Gateway” (G-3) mean that the following action is performed under the condition on the connected “Condition Flow”. The difference is that there is one and only one condition after the “Exclusive Gateway” can be met, while this does not apply to the “Inclusive Gateway”. The “Parallel Gateway” (G-5) represents that the following actions are performed in parallel. Note that, all gateways appear in pairs and the latter ones (G-4, G-2, G-6) indicate the end of the non-sequential activities. Additionally, the “Data Constraint” and “Action Constraint” represent the necessary data (C-1) and essential notices (C-2) for actions connected by the “Constraint Flow”, respectively.

3.2 Dataset Construction

With these high-quality procedural graphs, we then perform the dataset construction as a data2text task (Lin et al., 2023). Specifically, we design a three-stage pipeline: 1) Decomposition & Transformation that decomposes the graph into fragmented spans/sentences in natural language; 2) Grouping & Ordering that logically organizes the procedural fragments; 3) Aggregating & Smoothing that unifies the fragments into high-quality documents.

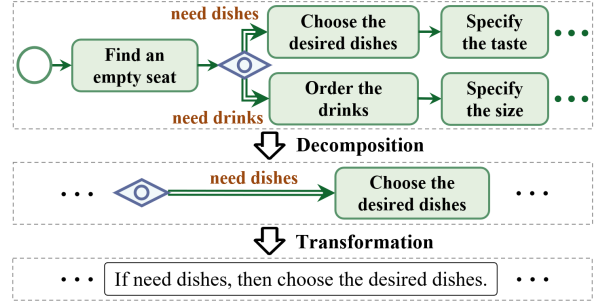


Figure 2: Illustration of decomposing the graph into units and transforming a unit into a procedural fragment.

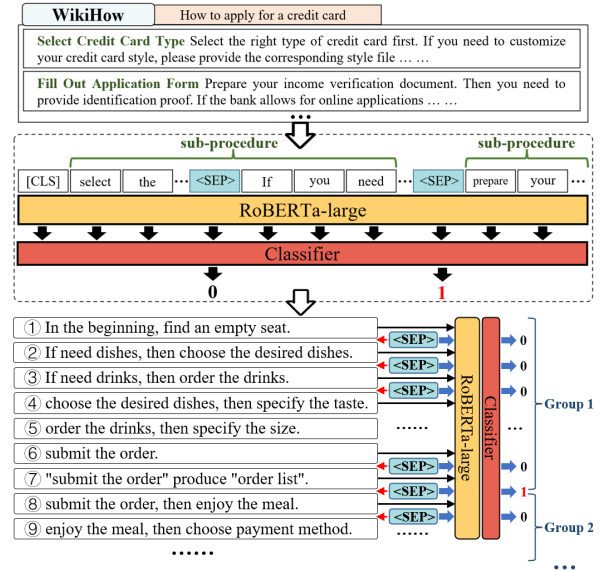


Figure 3: The grouping of procedural fragments using a pre-trained boundary identification model.

Decomposition & Transformation To narrow the huge gap between the complex graph and the length document, we first decompose the graph into minimal meaningful units. We define a vocabulary with nineteen units, each of which consists of actions, gateways or constraints connected by flow (cf., Appx. A). We also design nineteen hand-written templates for the vocabulary. Based on the unit vocabulary, we perform the breadth-first search strategy (Korf, 1985) over the graph, as shown in Figure 2. Particularly, some actions may be repeatedly walked to preserve the sequential execution relation between adjacent actions in the graph. For example, “Order drinks” is walked twice, forming “If need drinks, then order drinks.” and “Order drinks, then specify the size.”. We then transform the decomposed units into natural language spans/sentences based on the paired templates. We call these spans/sentences “procedural fragments”.

Grouping & Ordering When writing a procedural document, it is vital to provide information in a logical way, namely, describing the whole procedure sub-procedure by sub-procedure (Futrelle, 1999, 2004). Hence, we need to group the fragments into sub-procedures. Specifically, we employ a simple boundary identification model, which employs the RoBERTa-large model (Liu et al., 2019) with a classifier layer, to predict whether a fragment is the end of a group. We train it on the WikiHow corpus Bolotova et al. (2023), which consists of procedural documents collected from the WikiHow website and format marks indicate different sub-procedures. The pre-trained model is then used to assign group marks for fragments, cf., Figure 3. It is worth noticing that the way people describe a procedure is not the way the machine searches over the graph. Thus, it is necessary to order the fragments to better match human expression. For example, we need to exchange fragment ③ and fragment ④ of Group 1 in Figure 3. We achieve this via a fragment ordering model, whose structure is borrowed from Bin et al. (2023). Besides the WikiHow documents, we train this model on the remaining publicly available procedural document datasets (Castelli et al., 2020; Zhang et al., 2020; Lyu et al., 2021; Sakaguchi et al., 2021; Nandy et al., 2021) to reorder sentences after random shuffle. We use the pre-trained model to determine the order of fragments in the same group.

Aggregating & Smoothing In a standard procedural document, a single sentence (“If the customer needs dishes, then choose the desired dishes and specify the taste.”) may convey multiple procedural fragments (② and ④). Given this, we reuse the boundary identification model to aggregate fragments that should be presented in the same sentence. The model is retrained to identify the correct ends of sentences from randomly inserted ones on the datasets used in the ordering phase. The pre-trained model is used to assign sentence marks for the fragments. We further add the actor information before all fragments corresponding to each actor. At this point, all fragments have been organized in proper order with group and sentence marks. We paraphrase the fragments via ChatGPT, which has shown near or even superior human-level performance in many paraphrasing task (Chui, 2023). After paraphrasing, we notice that there are a small number of redundant expressions caused by repeatedly walked nodes and inconsistent ac-

tions/constraints generated by ChatGPT. Thus, we further refine the documents with a few handwriting rules and manual corrections. At last, we develop a dataset with 3,394 high-quality document-graph pairs. On average, each document contains 10.67 sentences and each sentence contains 15.22 words. See Table 4 in Appx. A for more statistics.

3.3 Dataset Analysis

We analyze our datasets to investigate whether the generated procedural documents are consistent with the original graphs, whether the documents are qualified according to human standards, and whether the proposed strategies contribute to a better quality of the dataset. Therefore, we conduct both automatic and human evaluations to compare the dataset constructed by our three-stage pipeline with the datasets constructed by two variations.

Automatic Evaluation We adopt two commonly used Data2Text metrics. The *FINE* score (Faille et al., 2021) models the evaluation as a natural language inference task — by inferring fragments from the documents it checks omissions, while the other direction checks hallucinations. The *ESA* score (Faille et al., 2021) evaluates the coverage of entities and actions in the documents. For both, higher values indicate better performance. Details are listed in the appendix B.

Human Evaluation Inspired by Miller (1979), we ask three workers to score the document from 1 to 5 on five criteria: *readability*, the quality of the document to be understood easily; *accuracy*, the quality of the document accurately describing the information in the procedure graph; *clarity*, the quality of the document expressing the complex execution of actions logically; *simplicity*, the quality of the document not containing redundant information; *usability*, the quality of the document aiding users in accomplishing this procedure. Details are listed in the appendix B.

Variations We create two variants of our three-stage pipeline: *concatenating*, which directly concatenates all fragments to form the documents; *paraphrasing*, which directly paraphrases the concatenation of all fragments using ChatGPT without grouping, ordering and aggregating process. Note that, we don’t involve *concatenating* in the automatic evaluation. This is because the concatenation of all fragments is always consistent with the

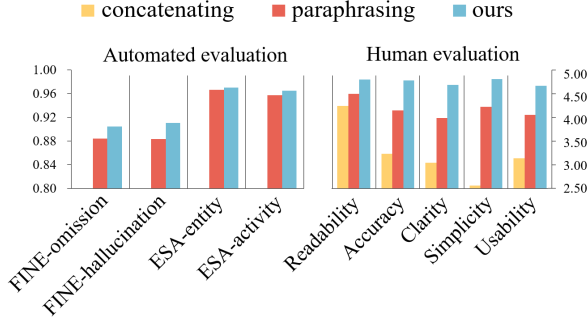


Figure 4: Comparison of our method with two variations via automatic and human evaluations.

information in the graphs but lacks fluency and logicity which requires further human evaluation.

Results As shown in Figure 4, the proposed three-stage pipeline achieves better performance than the other two variations under both automatic and human evaluation. Automatic evaluation reveals that although the paraphrasing variation obtains almost close ESA scores to our pipeline, it loses the game completely on the FINE metrics without the help of grouping, ordering and aggregating strategies. From human evaluation, we further demonstrate the superiority of our pipeline to describe procedural graphs in a fluent, accurate, logical, simple, and user-friendly way. Specifically, the concatenating variation gets the lowest scores on all criteria due to its inaction of unreadable spans, redundant information, and chaotic orders. In line with the automatic evaluation, we observe the largest gap between the paraphrasing variation and our pipeline on the clarity criterion due to its invisibility of complex logic among the fragments. Note that, the ICC scores of all evaluations are above 0.75, indicating the reliability of human results (Koo and Li, 2016).

4 Experiments

We conduct systematic experiments to answer *Q1* and *Q2* raised in Sec. 1. We split the dataset into train, validation and test sets with 3:1:2 ratio. For *Q1*, we collect state-of-the-art baselines and evaluate them in PAGED. We further introduce three metrics based on the underlying structure of graphs and the surface form of elements. Specifically, we use the BLEU scores to measure the performance of actor, action and constraint extraction and F1 scores to measure the performance of gateway prediction. Besides, the performance of flow prediction is measured via soft F1 scores (Tandon et al., 2020), which is computed based on the BLEU scores of

associated textual elements. For *Q2*, we involve Flan-T5 (Chung et al., 2022), ChatGPT (Ouyang et al., 2022) and Llama2 (Touvron et al., 2023) in PAGED to show their potentials and improve their performance using a self-refine strategy.

4.1 Performance of Baselines (*Q1*)

Baselines We collect five baselines:

- Sonbol et al. (2023) uses rules to extract sequential, exclusive, and parallel actions with a few data constraints.
- Neuberger et al. (2023) designs a pipeline to extract sequential actions and organize partial non-sequential actions, ignoring all constraints and inclusive gateways.
- Sholih et al. (2022) extracts actions, exclusive gateways, and parallel gateways based on a pre-defined representation of the procedural graph.
- PET (Bellan et al., 2023) trains a sequence tagging model to extract actions with a few constraints and uses rules to construct the final graph.
- CIS (Bellan et al., 2022) presents a rough attempt to use LLMs¹ for action extraction via few-shot in-context learning and constructs the graphs via handwritten rules.

Results As shown in Table 2 (Rows 1-5), existing studies are far from solving this task well, especially when organizing the logical structure of graphs (cf., the results on gateways and flows). This is because either rules or neural models are derived from limited data, leading to an incomplete coverage of all elements and an incomprehensive understanding of complex documents. Additionally, we have the following observations:

1) Heuristic methods (Row 1-3) perform poorly for actor, action and constraint extraction. The reason is that hand-written rules fail to understand various expressions and coreferences, resulting in poor generalization.

2) PET (Row 4), though being a customized deep neural model, only performs slightly better than heuristic methods. This is because the PET model is only trained on 45 samples and utilizes several rules to construct the flows.

3) We conjecture the reason why all baselines only meet part of the requirements of optimal pro-

¹The LLM used in Bellan et al. (2022) is GPT3. We replace it with ChatGPT in the current setting for a fair comparison.

Table 2: Performances of state-of-the-art baselines and LLMs. Higher values indicate better performances.

Row	Model	Actor	Action	Constraint		Gateway			Flow		
				Data	Action	Exclusive	Inclusive	Parallel	Sequence	Condition	Constraint
1	Sonbol et al. (2023)	0.028	0.308	0.213	-	0.485	-	0.279	0.056	0.047	0.017
2	Neuberger et al. (2023)	0.027	0.276	-	-	0.469	-	0.337	0.074	0.061	-
3	Sholiq et al. (2022)	-	0.387	-	-	0.463	-	0.198	0.091	0.022	-
4	PET (Bellan et al., 2023)	0.085	0.430	0.069	-	0.493	-	-	0.164	0.026	0.000
5	CIS (Bellan et al., 2022)	0.633	0.639	-	-	0.455	-	-	0.203	0.157	-
	Flan-T5 (Chung et al., 2022)										
6	+ Few-shot In-context Learning	0.206	0.362	-	-	0.376	-	-	0.084	0.013	-
7	+ Supervised Fine-tuning	0.659	0.684	0.589	<u>0.366</u>	0.419	0.045	<u>0.393</u>	0.395	<u>0.168</u>	0.363
8	ChatGPT (Ouyang et al., 2022)										
	+ Few-shot In-context Learning	0.625	0.681	<u>0.687</u>	0.286	0.477	0.173	0.388	<u>0.408</u>	0.158	<u>0.444</u>
	Llama2 (Touvron et al., 2023)										
9	+ Few-shot In-context Learning	0.502	0.573	0.357	0.049	0.446	0.067	0.128	0.193	0.107	0.201
10	+ Supervised Fine-tuning	0.674	0.744	0.779	0.499	0.554	<u>0.090</u>	0.398	0.478	0.319	0.467

* The best results are marked in bold, and the second-best results are marked with underlines.

cedural graphs lies in the huge cost of writing rules and annotating data. We believe this issue can be alleviated with the dataset created in Sec. 3.

4) It is not surprising that CIS (Row 5) achieves the highest scores on the action and actor extraction among these baselines, as it gains more power to understand word meanings from the LLM. Note that, the increase in flow prediction between PET and CIS also comes from more accurate action extraction instead of more delicate rules to construct the graph. This increase further encourages us to investigate more potentials of LLMs in procedural graph extraction.

5) All baselines perform poorly on flow prediction, indicating the challenge of understanding logical structures in documents. This motivates us to find out, besides an in-depth study of LLMs, what else matters to construct the logic in graphs.

4.2 Performance of LLMs (Q2)

LLMs We investigate three advanced LLMs. The first one is ChatGPT (Ouyang et al., 2022), which is good at information comprehension and text generation. Different from CIS (Bellan et al., 2022), we use our high-quality dataset to apply few-shot in-context learning (ICL) on ChatGPT. The other two (Flan-T5 (Chung et al., 2022) and Llama2 (Touvron et al., 2023)) are open-source alternatives to ChatGPT. Besides ICL, we also deploy the alternatives with supervised fine-tuning strategies.

Results As expected, the LLMs (Rows 6-10) update state-of-the-art results on all metrics, especially for actor, action and constraint extraction. However, for gateway and flow predictions, all LLMs can hardly get > 0.5 F1 scores, exhibiting their weakness in arranging logical structures of graphs. We also have the following observations:

1) Our high-quality dataset significantly promotes LLMs’ ability on procedural graph extraction. A piece of direct evidence is that we see a rising trend between the results on Rows 5 and 8, whose only difference lies in the data used in the few-shot ICL. Another piece of evidence is that the gap between Flan-T5 and Llama2 is rapidly narrowed after using more data from our dataset for tuning.

2) The procedural knowledge, especially the logical structure of non-sequential actions, has been overlooked during the initial training of LLMs. This is demonstrated by the fact that fine-tuned LLMs perform better than few-shot LLMs. The Llama2 model even beats ChatGPT after learning more procedural knowledge from supervised fine-tuning.

3) LLMs, including CIS, show significant potential for actor, action and constraint extraction, indicated by the large improvements compared with Rows 1-4. This is because LLMs are good at understanding lengthy contexts and thus have the advantage of identifying meaningful elements from procedural documents. We also believe that the emergence of more powerful LLMs in the future will continue to promote better results on these metrics.

4) We believe the biggest challenge for LLMs to extract accurate procedural graphs lies in the lack of logic reasoning ability among actions, especially non-sequentially executed ones. The supporting evidence is that, though largely surpassing baselines on the actor, action and constraint extractions, LLMs don’t present such impressive improvements in gateway extraction. This implies that we can boost LLMs’ performances by paying extra attention to non-sequential logic prediction.

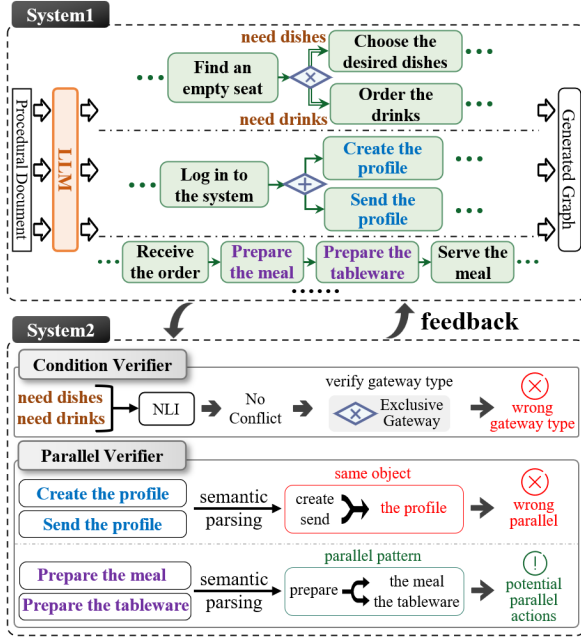


Figure 5: The self-refine strategy, in which “System1” extracts procedural graphs and “System2” verifies gateways of graphs and provides feedback for refinement.

4.3 Self-refine Strategy

To overcome the above-mentioned challenge, inspired by Nye et al. (2021); Madaan et al. (2023), we design a self-refine strategy to help LLMs gain logic reasoning ability among actions from iterative feedback and refinement. As shown in Figure 5, it consists of two systems — “System1” is used to directly extract procedural graphs from documents, “System2” is used to verify the extracted graphs and provide feedback for further refinement of “System1”. In “System2”, we center on the shortest slab and carefully examine the gateway prediction results using the condition and parallel verifiers.

Condition Verifier We design the condition verifier to handle both exclusive and inclusive gateways, whose key difference lies in the conditions followed by gateways. It is worth noticing that, with the exclusive gateway, there is always one and only one condition that can be met. Accordingly, we suppose if the conditions hold conflict, the gateway should be the exclusive one; otherwise, it should be the inclusive one. Particularly, we use a pre-trained natural language inference (NLI) model (Liu et al., 2019) to detect the conflict and verify gateways. For example, after feeding “need dishes” and “need drinks” into the NLI model, we get “No Conflict”. This suggests the gateway should be the inclusive one, which is different from the result predicted by “System1”. In this case, the

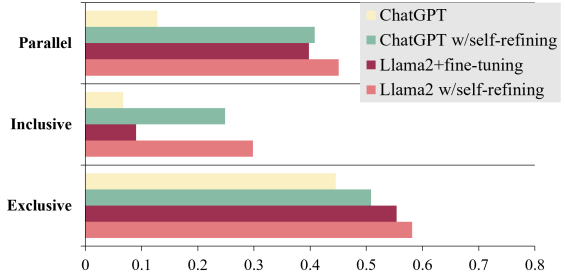


Figure 6: F1 scores on gateway predictions of LLMs and their variants with our self-refine strategy.

condition verifier is triggered to provide feedback to refine “System1”.

Parallel Verifier We further design the parallel verifier to reorganize the actions performed in parallel. We notice that if the actions act on the same object, they can never be executed in parallel. Given that, we extract the objects of actions using the semantic parsing tool and determine the parallel gateway based on these objects. For example, as “create the profile” and “send the profile” have the same object “the profile”, they cannot be performed in parallel. Besides, the optimal procedural graph is expected to help users in a time-saving way (Miller, 1979). Thus, we further prepare another type of feedback by examining the sequential actions. For example, as “prepare the meal” and “prepare the tableware” have different objects, they have a big chance to be performed in parallel. We provide “System1” with both types of feedback for the refinement of mistakes on parallel gateways.

Results We apply the self-refine strategy to the top-two winners in Table 2, i.e., ChatGPT and fine-tuned Llama2. As shown in Figure 6, both ChatGPT and fine-tuned Llama2 have better performances with the help of our self-refine strategy. More surprisingly, there are significant improvements in inclusive gateway extraction, whose previous scores are extremely poor. This indicates that, with effective strategies, LLMs have the potential to gain logic reasoning ability among actions including non-sequential ones. The improvement of Llama2 is not as large as that of ChatGPT. This is because the model’s gain from the self-refine strategy drops as it encounters more procedural knowledge. This suggests the importance of incorporating the learning of procedural knowledge during the pre-training stage of LLMs, which will be beneficial for LLMs’ logic reasoning ability.

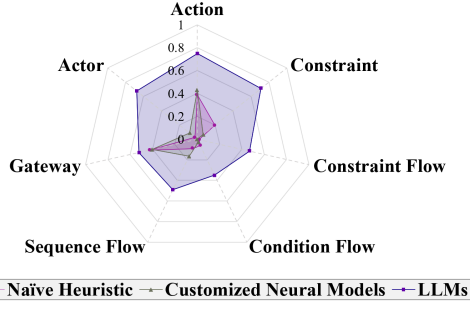


Figure 7: Best performances of heuristics, customized neural models and LLMs on seven dimensions.

Auxiliary Analysis We group the methods into three sets, i.e., the naïve heuristics (rows 1-3), customized neural models (row 4), and LLMs (rows 5-10). Towards a clear understanding of the advantages and challenges, we report the best performances of each type of method on seven dimensions in Figure 7. Even with the highest scores, heuristic methods and customized neural models exhibit poor performance across all dimensions and can hardly handle condition and constraint flows due to their neglect of the logical structure in documents. LLMs show substantial improvements compared to the others in all dimensions except for the gateway. This suggests that even the powerful LLMs face challenges in managing non-sequential actions, which is also the main challenge when conducting optimal procedural graph extraction.

5 Conclusion and Future work

We propose the PAGED benchmark, where we systematically study the progress of current procedural graph extraction methods and explore the potential of emerging LLMs on this task. We equip PAGED with a high-quality dataset, which is about ten times larger than the previous largest ones. Experimental results of baselines in PAGED reveal that current methods are far from solving this task well. We further involve three advanced LLMs in PAGED to demonstrate their advantages in extracting textual elements and challenges in organizing local structures. To overcome the main challenge, we design a novel self-refine strategy to empower the LLMs’ ability in reasoning gateways. The results show that, with effective strategies, LLMs have the potential of LLMs to comprehend the logical structure among non-sequential actions.

We hope PAGED can benefit the research on optimal procedural graphs extraction. There are several directions for further work. First, to improve the

performance of LLMs, we suggest introducing procedural knowledge during the pre-training stage of LLMs. Second, to gain better flow extraction, we will explore more effective methods for handling complex logic structures in the documents. Lastly, we plan to find a real-world scenario to investigate, besides accuracy, what else is limiting the practical usage of automatic procedure graph extraction.

6 Ethics Statement

The dataset we constructed is sourced from a publicly available model collection originated from the BPM Academic Initiative and does not contain any sensitive or personal privacy related information. The procedural graphs used in our dataset are available for research purposes on the “CC BY-NC-SA 3.0 DEED” licence², which explicitly permits that we can not only use the collected examples but also “remix, transform, and build upon the material”. Therefore, we believe that there is no ethical issue with our work.

7 Limitations

The procedural documents in our dataset can hardly be equal to the documents directly written by experts. Despite this fact, we believe our dataset can still largely promote the study of automatic procedural graph extraction. The samples for few-shot learning are randomly selected from the dataset. This may lead to fluctuations in LLMs’ results. However, since the samples of all LLMs are randomly chosen, we believe the experimental setup is fair. We acknowledge that carefully selecting samples would yield better results, but this is not the focus of this benchmark. We only design one prompt for all LLMs. Although using another elaborate prompt could introduce new variations to the experimental results, we consider this as a topic for future research.

Acknowledgement

This work was supported in part by the National Natural Science Foundation of China (No. 62206191 and No. 62272330); in part by the Natural Science Foundation of Sichuan (No. 2023NS-FSC0473), and in part by the Fundamental Research Funds for the Central Universities (No. 2023SCU12089 and No. YJ202219).

²<http://fundamentals-of-bpm.org/process-model-collections/>

References

- Lars Ackermann, Julian Neuberger, and Stefan Jablon-ski. 2021. Data-driven annotation of textual process descriptions based on formal meaning representations. In *International Conference on Advanced Information Systems Engineering*, pages 75–90. Springer.
- Patrizio Bellan, Mauro Dragoni, and Chiara Ghidini. 2022. Leveraging pre-trained language models for conversational information seeking from text. *arXiv e-prints*, pages arXiv–2204.
- Patrizio Bellan, Han van der Aa, Mauro Dragoni, Chiara Ghidini, and Simone Paolo Ponzetto. 2023. Pet: An annotated dataset for process extraction from natural language text tasks. *LECTURE NOTES IN BUSINESS INFORMATION PROCESSING*, 460:315–321.
- Yi Bin, Wenhao Shi, Bin Ji, Jipeng Zhang, Yujuan Ding, and Yang Yang. 2023. Non-autoregressive sentence ordering. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4198–4214.
- Valeriia Bolotova, Vladislav Blinov, Sofya Filippova, Falk Scholer, and Mark Sanderson. 2023. Wikihowqa: A comprehensive benchmark for multi-document non-factoid question answering. In *Proceedings of the 61th Conference of the Association for Computational Linguistics*.
- Vittorio Castelli, Rishav Chakravarti, Saswati Dana, Anthony Ferritto, Radu Florian, Martin Franz, Dinesh Garg, Dinesh Khandelwal, J Scott McCarley, Michael McCawley, et al. 2020. The techqa dataset. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1269–1278.
- Ho Chui Chui. 2023. Chatgpt as a tool for developing paraphrasing skills among esl learners. *Journal of Creative Practices in Language Learning and Teaching (CPLT)*, 11(2).
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Marlon Dumas, Marcello La Rosa, Jan Mendling, Hajo A Reijers, et al. 2018. *Fundamentals of business process management*, volume 2. Springer.
- Song Duong, Alberto Lumbreras, Mike Gartrell, and Patrick Gallinari. 2023. Learning from multiple sources for data-to-text and text-to-data. In *International Conference on Artificial Intelligence and Statistics*, pages 3733–3753. PMLR.
- Ondřej Dušek and Zdeněk Kasner. 2020. Evaluating semantic accuracy of data-to-text generation with natural language inference. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 131–137.
- Elena Viorica Epure, Patricia Martín-Rodilla, Charlotte Hug, Rebecca Deneckère, and Camille Salinesi. 2015. Automatic process model discovery from textual methodologies. In *2015 IEEE 9th International Conference on Research Challenges in Information Science (RCIS)*, pages 19–30. IEEE.
- Juliette Faille, Albert Gatt, and Claire Gardent. 2021. Entity-based semantic adequacy for data-to-text generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1530–1540.
- Renato César Borges Ferreira, Lucinéia Heloisa Thom, and Marcelo Fantinato. 2017. A semi-automatic approach to identify business process elements in natural language texts. In *International Conference on Enterprise Information Systems*, volume 2, pages 250–261. SCITEPRESS.
- Fabian Friedrich, Jan Mendling, and Frank Puhlmann. 2011. Process model generation from natural language text. In *Advanced Information Systems Engineering: 23rd International Conference, CAiSE 2011, London, UK, June 20-24, 2011. Proceedings 23*, pages 482–496. Springer.
- Zihao Fu, Bei Shi, Wai Lam, Lidong Bing, and Zhiyuan Liu. 2020. Partially-aligned data-to-text generation with distant supervision. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9183–9193.
- Robert P Futrelle. 1999. Summarization of diagrams in documents. *Advances in Automated Text Summarization*, pages 403–421.
- Robert P Futrelle. 2004. Handling figures in document summarization. In *Text Summarization Branches Out*, pages 61–65.
- Emden Gansner, Eleftherios Koutsofios, and Stephen North. 2006. Drawing graphs with dot.
- Joachim Herbst and D Karagiannis. 1999. An inductive approach to the acquisition and adaptation of workflow models. In *Proceedings of the IJCAI*, volume 99, pages 52–57. Citeseer.
- Krzysztof Honkisz, Krzysztof Kluza, and Piotr Wiśniewski. 2018. A concept for generating business process models from natural language description. In *Knowledge Science, Engineering and Management: 11th International Conference, KSEM 2018, Changchun, China, August 17–19, 2018, Proceedings, Part I 11*, pages 91–103. Springer.
- 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*.
- Zdeněk Kasner and Ondřej Dušek. 2022. Neural pipeline for zero-shot data-to-text generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3914–3932.

- Terry K Koo and Mae Y Li. 2016. A guideline of selecting and reporting intraclass correlation coefficients for reliability research. *Journal of chiropractic medicine*, 15(2):155–163.
- Richard E. Korf. 1985. [Depth-first iterative-deepening: an optimal admissible tree search](#). *Artif. Intell.*, 27(1):97–109.
- Hongru Liang, Jia Liu, Weihong Du, Dingnan Jin, Wenqiang Lei, Zujie Wen, and Jiancheng Lv. 2023. Knowing-how & knowing-that: A new task for machine comprehension of user manuals. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 10550–10564.
- Yupian Lin, Tong Ruan, Jingping Liu, and Haofen Wang. 2023. A survey on neural data-to-text generation. *IEEE Transactions on Knowledge and Data Engineering*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Hugo A López, Rasmus Strømsted, Jean-Marie Niyodusenga, and Morten Marquard. 2021. Declarative process discovery: Linking process and textual views. In *International Conference on Advanced Information Systems Engineering*, pages 109–117. Springer.
- Qing Lyu, Li Zhang, and Chris Callison-Burch. 2021. Goal-oriented script construction. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 184–200.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651*.
- Bilal Maqbool, Farooque Azam, Muhammad Waseem Anwar, Wasi Haider Butt, Jahan Zeb, Iqra Zafar, Aiman Khan Nazir, and Zuneera Umair. 2019. A comprehensive investigation of bpmn models generation from textual requirements—techniques, tools and trends. In *Information Science and Applications 2018: ICISA 2018*, pages 543–557. Springer.
- Jan Mendling, Henrik Leopold, Lucineia Heloisa Thom, and Han van der Aa. 2019. Natural language processing with process models (nlp4re report paper). In *REFSQ Workshops*.
- Carolyn R Miller. 1979. A humanistic rationale for technical writing. *College English*, 40(6):610–617.
- Yoshio Momouchi. 1980. Control structures for actions in procedural texts and pt-chart. In *COLING 1980 Volume 1: The 8th International Conference on Computational Linguistics*.
- Abhilash Nandy, Soumya Sharma, Shubham Madhaskiya, Kapil Sachdeva, Pawan Goyal, and Niloy Ganguly. 2021. Question answering over electronic devices: A new benchmark dataset and a multi-task learning based qa framework. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4600–4609.
- Julian Neuberger, Lars Ackermann, and Stefan Jablon-ski. 2023. Beyond rule-based named entity recognition and relation extraction for process model generation from natural language text. *arXiv preprint arXiv:2305.03960*.
- Maxwell Nye, Michael Tessler, Josh Tenenbaum, and Brenden M Lake. 2021. Improving coherence and consistency in neural sequence models with dual-system, neuro-symbolic reasoning. *Advances in Neural Information Processing Systems*, 34:25192–25204.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Kuntal Kumar Pal, Kazuaki Kashiara, Pratyay Banerjee, Swaroop Mishra, Ruoyu Wang, and Chitta Baral. 2021. Constructing flow graphs from procedural cybersecurity texts. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3945–3957.
- Chen Qian, Lijie Wen, Akhil Kumar, Leilei Lin, Li Lin, Zan Zong, Shu’ang Li, and Jianmin Wang. 2020. An approach for process model extraction by multi-grained text classification. In *Advanced Information Systems Engineering: 32nd International Conference, CAiSE 2020, Grenoble, France, June 8–12, 2020, Proceedings 32*, pages 268–282. Springer.
- Luis Quishpi, Josep Carmona, and Lluís Padró. 2020. Extracting annotations from textual descriptions of processes. In *International Conference on Business Process Management*, pages 184–201.
- Sebastian Raschka. 2018. Model evaluation, model selection, and algorithm selection in machine learning. *arXiv preprint arXiv:1811.12808*.
- Haopeng Ren, Yushi Zeng, Yi Cai, Bihan Zhou, and Zetao Lian. 2023. Constructing procedural graphs with multiple dependency relations: A new dataset and baseline. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8474–8486.
- Keisuke Sakaguchi, Chandra Bhagavatula, Ronan Le Bras, Niket Tandon, Peter Clark, and Yejin Choi. 2021. proscripT: Partially ordered scripts generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2138–2149.
- Sholiq Sholiq, Riyanarto Sarno, and Endang Siti As-tuti. 2022. Generating bpmn diagram from textual

- requirements. *Journal of King Saud University-Computer and Information Sciences*, 34(10):10079–10093.
- Patrick E Shrout and Joseph L Fleiss. 1979. Intraclass correlations: uses in assessing rater reliability. *Psychological bulletin*, 86(2):420.
- Riad Sonbol, Ghaida Rebdawi, and Nada Ghneim. 2023. A machine translation like approach to generate business process model from textual description. *SN Computer Science*, 4(3):291.
- Yixuan Su, Zaiqiao Meng, Simon Baker, and Nigel Collier. 2021. Few-shot table-to-text generation with prototype memory. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 910–917.
- Niket Tandon, Keisuke Sakaguchi, Bhavana Dalvi, Dheeraj Rajagopal, Peter Clark, Michal Guerquin, Kyle Richardson, and Eduard Hovy. 2020. [A dataset for tracking entities in open domain procedural text](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6408–6417, Online. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Mark von Rosing, Stephen White, Fred Cummins, and Henk de Man. 2015. Business process model and notation-bpmn.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Rong Ye, Wenxian Shi, Hao Zhou, Zhongyu Wei, and Lei Li. 2019. Variational template machine for data-to-text generation. In *International Conference on Learning Representations*.
- Li Zhang, Qing Lyu, and Chris Callison-Burch. 2020. [Reasoning about goals, steps, and temporal ordering with WikiHow](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4630–4639, Online. Association for Computational Linguistics.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.

A Details of Data Construction

Decomposition & Transformation To narrow the huge gap between the complex graph and the length document and meanwhile maintain the consistency between the transferred fragments and the original graphs, such as the texts of the entities, actions, etc., we design hand-written templates to transfer the graphs into natural language fragments. The designed templates are listed in Table 3. Additionally, we find that not all units on the original graphs are meaningful and need to be explicitly expressed in generated procedural documents. So we filter out those meaningless units through heuristic rules.

Aggregating & Smoothing We conduct rephrasing operation on the concatenation of the processed fragments using LLM. We add the unique separator token “<SEP>” to indicate the group and sentence marks of the fragments, which can prompt the model to describe these fragments logically and coherently. Moreover, we add the actor information before all fragments corresponding to each actor. For example, we add the text “For the customer:” before all fragments corresponding to the customer. We adopt ChatGPT (gpt-3.5-turbo) to conduct the rephrasing operation. The designed prompt is shown in Figure 9. At last, we develop a dataset with 3,394 high-quality document-graph pairs. Additionally, although it is difficult to exactly define the large of the dataset, we argue our dataset is large enough because it is about ten times larger than the previous largest datasets and successfully supports the evaluation of current studies and the discovery of future directions in this field. We present the statistics information of the constructed dataset in Table 4.

B Details of Dataset Quality Evaluation

Automatic Evaluation Following the evaluation strategies commonly used in Data2Text task (Lin et al., 2023), we conduct automatic evaluation to evaluate whether the generated procedural documents provide information consistent with the original graphs. The automatic metrics used for evaluation are listed as follows:

FINE (Dušek and Kasner, 2020): evaluating the semantic equivalence of generated documents with a natural language inference model³. The natural

language inference model can be used to determine whether a “hypothesis” is true (entailment) given a “premise”. We first compute the entailment between the generated document and the transferred fragments to evaluate whether the generated document fully covers the original graph’s information (omission). We take the generated document as the “premise” and use the natural language inference model to determine whether the transferred fragments are true (entailment), i.e., whether the semantics of the transferred fragments are covered by the generated document. Then we exchange their positions to evaluate whether the generated document not contains redundant information beyond the graph (hallucination).

ESA (Faille et al., 2021): evaluating the faithfulness of the generated documents based on the coverage of the entities and actions in the original graphs. We use named entity recognition tool⁴ to extract the entities in the original graphs. Then we evaluate whether the entities and actions in the original graph are covered by the generated document with exact lexical match.

Human Evaluation We design five criteria to evaluate the generated documents through human scoring. To ensure the reliability of human evaluation, we hire three domain experts to evaluate the generated documents and calculate the ICC (Intraclass Correlation Coefficient) (Shrout and Fleiss, 1979) score between different experts. The higher ICC score indicates higher consistency between different experts and higher reliability of the evaluation. Generally an ICC of 0.75 or higher indicates that the evaluation is reliable (Koo and Li, 2016). ICC is calculated as follows:

$$ICC = \frac{MS_{\text{between}} - MS_{\text{within}}}{MS_{\text{between}} + (k - 1) \times MS_{\text{within}}} \quad (1)$$

where MS_{between} is the mean square for between groups variability, MS_{within} is the mean square for within groups variability and k is the number of groups. And we ensure that the ICC scores for all of our evaluations are ≥ 0.75 .

³<https://huggingface.co/FacebookAI/>

⁴[roberta-large-mnli](https://huggingface.co/roberta-large-mnli)

⁴<https://huggingface.co/dslim/bert-base-NER>

Table 3: Designed templates used to transfer the decomposed units into natural language procedural fragments. Due to the fact that the gateways in the procedural graphs are paired, each pair of gateways includes a “branch gateway” representing the beginning of the non-sequential execution of actions and a “merge gateway” representing the end of the non-sequential execution. We use prefix “B_” to indicate the branch gateway, and prefix “M_” to indicate the merge gateway. In addition, the “Flow” in a unit is used to connect different elements, while the “condition” in a unit represents that this Flow connects exclusive or inclusive gateways with other elements, so there exists specific condition on this Flow. For simplicity, we use “XOR”, “OR” and “AND” as abbreviations for exclusive, inclusive and parallel gateways respectively.

Decomposed Unit	Template
(Start, Flow, Action)	In the beginning, {Action}.
(Start, Flow, B_Gateway)	In the beginning,
(Action1, Flow, Action2)	{Action1}, then {Action2}.
(Action, Flow, End)	{Action}, and the procedure ends.
(B_XOR, condition, Action)	If {condition}, then {Action}.
(B_XOR, condition, B_Gateway)	If {condition},
(B_XOR, condition, M_Gateway, Flow, Action)	If {condition}, then {Action}.
(B_XOR, condition, End)	If {condition}, then the procedure ends.
(B_OR, condition, Action)	If {condition}, then {Action}.
(B_OR, condition, B_Gateway)	If {condition},
(B_OR, condition, M_Gateway, Flow, Action)	If {condition}, then {Action}.
(B_OR, condition, End)	If {condition}, then the procedure ends.
(B_AND, Flow, * ¹) (B_AND, Flow, * ²)	{* ¹ }, at the same time, {* ² }.
(B_AND, Flow, * ¹) (B_AND, Flow, * ²) (B_AND, Flow, * ³)	{* ¹ }, at the same time, {* ² }, meanwhile, {* ³ }.
(M_Gateway, Flow, Action)	{Action}.
(M_Gateway, Flow, End)	The procedure ends.
(Action, Flow, DataConstraint)	“{Action}” produce “{DataConstraint}”.
(DataConstraint, Flow, Action)	“{Action}” require access to “{DataConstraint}”.
(Action, Flow, ActionConstraint)	For {Action}, pay attention to that {ActionConstraint}.

Table 4: Statistics of the constructed dataset. We present statistical information on the number of documents and various elements in the dataset.

Statistics	Num	Statistics	Num
Document	3394	Data Constraint	3500
Sentence	37226	Action Constraint	2307
Token	566639	Sequence Flow	36438
Action	36537	Condition Flow	10598
Exclusive Gateway	7024	Constraint Flow	5807
Inclusive Gateway	1204	Actor	22775
Parallel Gateway	2050		

C Details of Experiments

C.1 Model Details

Heuristic models do not require training data, so we conduct evaluation for heuristic models according to the original papers’ setting. The implementation details of other models are listed as follows:

- **PET**: We use Roberta-large (Liu et al., 2019) as the backbone-model and train the model on

the publicly available dataset ⁵. A total of 3 epochs are trained, the adopted optimizer is AdamW and the learning rate is set to 5e-6.

- **CIS**: We use ChatGPT (Ouyang et al., 2022) as the pre-trained language representation model for in-context learning. We use the best-performed prompt templates following the original paper to conduct extraction with the model.
- **Flan-T5**: We fine-tune the Flan-t5-xxl model on our train set data with low-rank adaptation strategy Lora (Hu et al., 2021). A total of 10 epochs are trained and the learning rate is set to 1e-4.
- **Llama2**: Similar to Flan-T5 model, we fine-tune the Llama-2-70b-chat-hf model on our train set data with Lora. A total of 10 epochs are trained and the learning rate is set to 2e-4.

⁵<https://huggingface.co/datasets/patriziobellani/PET>

- **ChatGPT:** We use the gpt-3.5-turbo model to conduct conversation through the API provided by OpenAI⁶. We set temperature to zero and fix seed as 42 to eliminate the influence of random sampling of the model and enable stable reproduction.

C.2 Implementation Details

We split the dataset into train, validation and test sets with 3:1:2 ratio and evaluate the performance of all models uniformly on the test set. Heuristic models do not require training data, and PET trains the model on the sequence tagging data provided by itself. Therefore, we directly conduct evaluation for these models on our test set according to the original papers’ setting. End to End models utilize our train set data to fine-tune the model or construct examples for few-shot in-context learning, and use the validation set for model selection (Raschka, 2018). All the training process are conducted on a machine with $8 \times$ NVIDIA RTX A6000 GPUs. We use Copilot as an aid for coding.

To facilitate the LLMs to extract the procedural graphs with text generation, we use a dot language (Gansner et al., 2006) based graph representation in the form of “Element -> (condition) Element” to represent the extracted graphs. We use “XOR”, “OR” and “AND” as abbreviations for exclusive, inclusive and parallel gateways respectively, and use numbers as suffixes to distinguish different gateways of the same type on the graph (e.g., “OR1 -> (need dishes) Choose the desired dishes”). Additionally, we require the model to output the actors of corresponding extracted actions for actors predictions. Moreover, to make good use of the capabilities of LLMs, we adopt the few-shot in-context learning (ICL) and Chain-of-thought (CoT) (Wei et al., 2022) strategy to guide the reasoning process of extracting procedural graphs from documents for the LLMs, especially for the organization of non-sequential actions. The adopted prompt consisting of elaborate instruction and three examples is shown in Figure 10.

C.3 Metrics Details

We introduce three metrics based on the underlying structure of graphs and the surface form of elements to evaluate the model’s extraction of procedural graphs.

For actor, action and constraint extraction, we adopt the F1 based BLEU score (Liang et al., 2023) to evaluate how accurately can the model extract the texts of these three types of textual elements from the documents. Specifically, we first compute the BLEU score for each extracted actor, action or constraint by the model according to the elements with the same type in the gold procedural graph (e.g., for an extracted action, we find the most similar action in the gold procedural graph and compute the extracted action’s BLEU score based on the most similar action), thus calculating the precision scores of these extracted three types of textual elements. Then we do the same computation for all actors, actions and constraints in the gold procedural graph according to the extracted elements by the model to calculate the recall scores of all these three types of textual elements in the gold procedural graph. Then we can calculate the F1 scores based on the precision scores and recall scores. Note that we calculate the F1 scores of actors based on the best matched action pair (e.g., we first find the most similar action in the gold procedural graph for an extracted action and then calculate the BLEU score of the extracted actor by comparing these two actions’ actors).

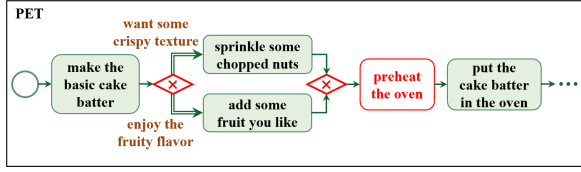
For gateway extraction, we adopt the standard F1 scores to evaluate whether the model can correctly use gateways to organize the non-sequential actions in the procedural graphs. Due to the fact that the gateways are meaningful only when paired with the corresponding elements in the graphs (von Rosing et al., 2015; Dumas et al., 2018), we consider an extracted gateway is correct only if its type and at least one of its paired element match those of the gold procedural graph (e.g., an extracted exclusive gateway and an action connecting with it by flow match a pair of such exclusive gateway and action in the gold procedural graph, this extracted exclusive gateway is considered correctly extracted). And for soft evaluation of the gateway extraction, an action or constraint is considered matched with the gold procedural graph if its BLEU score ≥ 0.5 .

For flow extraction, we measure the performance via a soft metric (Tandon et al., 2020), which computes the F1 scores based on the BLEU scores of associated textual elements. Specifically, a flow is considered correctly extracted only if its type and connected elements match those of the gold procedural graph. For condition flow, we compute the BLEU score between its condition and the gold

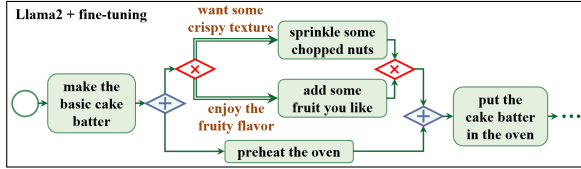
⁶<https://platform.openai.com/docs/api-reference/chat/create>

To make a cake by yourself, firstly, make the basic cake batter. If want some crispy texture, sprinkle some chopped nuts. If you enjoy the fruity flavor, add some fruit you like. Don't forget to preheat the oven while preparing these ingredients. After that, put the cake batter in the oven and ...

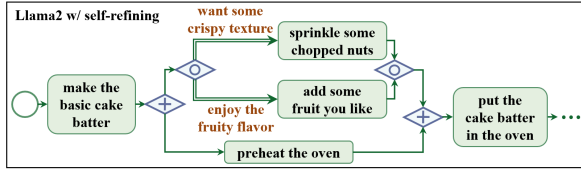
(a) The document to be extracted



(b) Screenshot of the procedural graph extracted by PET



(c) Screenshot of the procedural graph extracted by fine-tuned Llama2



(d) Screenshot of the procedural graph extracted by fine-tuned Llama2 with self-refine strategy

Figure 8: Illustration of case study.

label’s condition if another condition flow in the gold procedural graph can be matched.

C.4 Details of Improvement Strategies

Condition Verifier The designed templates used to provide feedback to the model are shown in Figure 10(a) and Figure 10(b). We feed the generated feedback based on the designed templates and extraction results to the model, so that the model can refine its extraction according to our provided feedback.

Parallel Verifier We use the off-the-shelf parsing tool to parse the syntactic structure of the text and obtain the predicate and object corresponding to each action. The designed templates used to provide feedback to the model are shown in Figure 10(c) and Figure 10(d).

C.5 Case Study

As shown in Figure 8, PET 8(b) fails to handle the parallel actions in the document 8(a) due to the lack of the ability to understand complex expressions, and it uses the wrong type of gateway

to organize the non-sequential actions “sprinkle some chopped nuts” and “add some fruit you like” as it can only deal with partial types of gateways. Fine-tuned Llama2 8(c) successfully organizes the parallel execution of actions in the document, but also uses the wrong gateway to organize the non-sequential actions “sprinkle some chopped nuts” and “add some fruit you like”. With the help of our designed verifier 8(d), fine-tuned Llama2 with self-refine strategy correct the wrong gateway into inclusive gateway through the feedback provided by the verifier, as there exists no conflict between the conditions “want some crispy texture” and “enjoy the fruity flavor”. This demonstrates the effectiveness of our proposed strategy by paying extra attention to non-sequential logic prediction.


```

"""Convert the following paragraph into fluent natural language document without changing its meaning.
The separator token <SEP> indicates the boundary between two parts of texts and it should not appear in the
output.
Please do not change the entities or actions mentioned in the paragraph and do not add any new entities or actions.
Make sure that those entities and actions in the generated document are exactly the same as those in the original
paragraph. Just make the paragraph more fluent and natural.

###
Paragraph:
{the concatenation of the processed fragments with separator tokens}

###
Output:
"""

```

Figure 9: The designed prompt used to rephrase the concatenation of the processed fragments.

For the following inclusive gateway, there exists conflict between the conditions followed by the gateway, which may indicate that it is impossible to satisfy multiple conditions simultaneously:

{extracted graph snippet containing the gateway}

Please check the conflict between these conditions and correct the gateway type to exclusive gateway (XOR) if there exists conflict between them.

(a) Template for inclusive gateway refinement

For the following exclusive gateway, there is no obvious conflict between the conditions followed by the gateway, which may indicate that it is possible to satisfy multiple conditions simultaneously:

{extracted graph snippet containing the gateway}

Please check the conflict between these conditions and correct the gateway type to inclusive gateway (OR) if there exists no conflict between them.

(b) Template for exclusive gateway refinement

For the following parallel gateway, these actions share the same object, which may prevent their parallel execution:

{extracted graph snippet containing the gateway}

Please check the difference of the objects for these parallel actions and convert them into sequential actions if these actions are executed on the same object.

(c) Template for parallel gateway refinement

For the following sequential actions, there exists obvious differences between the objects of these actions, which may support parallel execution of these actions:

{extracted graph snippet containing these sequential actions}

Please check the difference of the objects for these sequential actions and convert them into parallel actions if they support parallel execution.

(d) Template for sequential actions refinement

Figure 10: Designed templates of our proposed verifiers.

```

"""I want you to generate the Procedural Graph based on a Procedural Document.
The Procedural Graph contains the following types of "Nodes" and "Flows":

"Nodes":
"Start": start node indicates the start of a procedure, represented as "Start".
"End": end node indicates the ending of a procedure, represented as "End".
"Action": action node indicates a specific step in a procedure, represented as the step itself, such as "prepare the ingredients".
"XOR": exclusive gateway, indicates that only one of the following non-sequential actions can be executed, distinguish by numbers, such as "XOR1".
"OR": inclusive gateway, indicates that one or more of the following non-sequential actions can be executed, distinguish by numbers, such as "OR1".
"AND": parallel gateway, indicates that all of the following actions should be executed in parallel, distinguish by numbers, such as "AND1".
"DataConstraint": DataConstraint indicates the constraints for the necessary data of the actions, represented as "DataConstraint(data object)".
"ActionConstraint": ActionConstraint indicates essential notices need to be considered for the execution of the actions, represented as "ActionConstraint(essential notices)".

" Flows":
"SequenceFlow": flow that represents the execution of sequential actions, such as "Start -> prepare the ingredients".
"ConditionFlow": the condition flow is used to indicate that the following action is performed under the condition on the Condition Flow, such as "XOR1 -> (condition1) choose the first one".
"ConstraintFlow": flow that is used to connect the constraints with corresponding actions, such as "prepare the ingredients -> ActionConstraint(essential notices)".

In addition, the actor of corresponding actions is put in the front of corresponding elements to indicate the actor of the following actions if needed, such as "For actor1.".

You should generate the graph in the format of "Node -> Node" line by line until generating the whole graph for the given Procedural Document, and keep the text of the nodes and conditions consistent with the original Procedural Document.

Here are some examples:

###
"Procedural Document":
Firstly, the customer needs to find an empty seat. If the customer needs dishes, then choose the desired dishes and specify the taste. If the customer needs drinks, then order the drinks and specify the size. The customer then submits the order, which is added to the order list. After enjoying the meal, the customer should choose the payment method. If the credit card is available, the customer pays by credit card; else if the credit card is not available, the customer should pay in cash. For the restaurant, once receiving the order from the order list, it prepares the meal according to the order and prepares the tableware for the customer at the same time. The meal is then served for the customer to enjoy. After that, the restaurant asks the customer to pay for the order and then confirms the payment. Note that the restaurant should provide the receipt if the customer needs. And the procedure ends.

###
"Procedural Graph":
The procedure starts with the customer finding an empty seat. Then there is an inclusive gateway to indicate the non-sequential actions because the customer can need dishes or drinks or both, and the customer should specify the taste after choosing the desired dishes and specify the size after ordering the drinks. Then the customer should submit the order to produce the order list data and enjoy the meal and then choose payment method. Then there is an exclusive gateway to indicate the non-sequential actions because the credit card is available or not. If the credit card is available, the customer pays by credit card; else if the credit card is not available, the customer should pay in cash. And the procedure for the customer ends. Then for the restaurant, there is a parallel gateway after receiving the order from the order list to indicate the non-sequential actions because the restaurant should prepare the meal and prepare the tableware in parallel. Then the restaurant serves the meal, asks the customer to pay for the order and confirms the payment. And note that provide the receipt if the customer needs when confirming the payment. And the procedure for the restaurant ends.
So the Procedural Graph of this Procedural Document is as follows:

For the customer:
Start -> find an empty seat
find an empty seat -> OR1
OR1 -> (needs dishes) choose the desired dishes
OR1 -> (needs drinks) order the drinks
choose the desired dishes -> specify the taste
order the drinks -> specify the size
specify the taste -> OR2
specify the size -> OR2
OR2 -> submits the order
submits the order -> DataConstraint(order list)
submits the order -> enjoy the meal
enjoy the meal -> choose payment method
choose payment method -> XOR1
XOR1 -> (credit card is available) pay by credit card
XOR1 -> (credit card is unavailable) pay in cash
pay by credit card -> XOR2
pay in cash -> XOR2
XOR2 -> End

For the restaurant:
Start -> receive an order
DataConstraint(order list) -> receive an order
receive an order -> AND1
AND1 -> prepare the meal
AND1 -> prepare the tableware
prepare the meal -> AND2
prepare the tableware -> AND2
AND2 -> serve the meal
serve the meal -> ask the customer to pay for the order
ask the customer to pay for the order -> confirm the payment
confirm the payment -> ActionConstraint(provide the receipt if the customer needs)
confirm the payment -> End

###
"Procedural Document":
In the beginning, the staff will receive an order request, and then checks the order type. If the order is standard type, the sufficiency of the stock is checked according to the stock table. If the order is special type, upload the order to the factory system. If the stock is sufficient for standard order, the goods will be directly shipped out, else if the stock is insufficient, they will need to be transferred from other warehouses. After that, the staff updates the order status and provide order information to the user. At the same time, the staff needs to bind order information to user account. Finally, the staff record the request status and the procedure ends.

###
"Procedural Graph":
The procedure starts with the staff receive an order request. After checking the order type, there is an exclusive gateway to indicate the non-sequential actions because the order can be standard type or special type. If the order is special type, upload the order to the factory system. And if the order is standard type, the sufficiency of the stock is checked according to the stock table data. And there is one more exclusive gateway after checking the sufficiency of the stock to indicate the non-sequential actions because the stock can be sufficient or insufficient. If the stock is sufficient for standard order, the goods will be directly shipped out, else if the stock is insufficient, they will need to be transferred from other warehouses. Then there is a parallel gateway to indicate the non-sequential actions because the staff should update the order status and provide order information to the user and meanwhile, bind order information to user account. Finally, the staff records the request status and the procedure ends.
So the Procedural Graph of this Procedural Document is as follows:

For the staff:
Start -> receive an order request
receive an order request -> check the order type
check the order type -> XOR1
XOR1 -> (the order is standard type) check the sufficiency of the stock
XOR1 -> (the order is special type) upload the order to the factory system
DataConstraint(the stock table) -> check the sufficiency of the stock
check the sufficiency of the stock -> XOR2
XOR2 -> (the stock is sufficient) directly shipped out the goods
XOR2 -> (the stock is insufficient) transfer the goods from other warehouses
directly shipped out the goods -> XOR3
transfer the goods from other warehouses -> XOR3
XOR3 -> XOR4
upload the order to the factory system -> XOR4
XOR4 -> AND1
AND1 -> update the order status
update the order status -> provide order information to the user
AND1 -> bind order information to user account
provide order information to the user -> AND2
bind order information to user account -> AND2
AND2 -> record the request status
record the request status -> End

###
"Procedural Document":
Start the service by receiving the email from the electronic mailbox, then parse the email content. If the email contains account query request, reply the account information to the user. If the email contains account modification request, record the information needs to be modified. After that, verify the validity of the account and verify the legality of the modified information at the same time if there exists account information to be modified. Otherwise update the verification timestamp of the account directly. Finally, synchronize the email content to the system and the procedure ends.

###
"Procedural Graph":
The procedure starts with receiving the email from the electronic mailbox data. Then parse the email content and there is an inclusive gateway to indicate the non-sequential actions because the email can contain account query request or account modification request or both. If the email contains account query request, reply the account information to the user. If the email contains account modification request, record the information needs to be modified, update the verification timestamp of the account directly. Else if there exists account information to be modified, there is a parallel gateway to indicate the non-sequential actions because we should verify the validity of the account and verify the legality of the modified information in parallel. Then synchronize the email content to the system and the procedure ends.
So the Procedural Graph of this Procedural Document is as follows:

Start -> receive the email
DataConstraint(electronic mailbox) -> receive the email
receive the email -> parse the email content
parse the email content -> OR1
OR1 -> (the email contains account query request) reply the account information to the user
OR1 -> (the email contains account modification request) record the information needs to be modified
reply the account information to the user -> OR2
record the information needs to be modified -> OR2
OR2 -> XOR1
XOR1 -> (there exists account information to be modified) AND1
XOR1 -> (otherwise) update the verification timestamp of the account directly
AND1 -> verify the validity of the account
AND1 -> verify the legality of the modified information
verify the validity of the account -> AND2
verify the legality of the modified information -> AND2
AND2 -> XOR2
update the verification timestamp of the account directly -> XOR2
XOR2 -> synchronize the email content to the system
synchronize the email content to the system -> End

Now you need to generate the corresponding Procedural Graph of the following Procedural Document:

###
"Procedural Document":
{}

###
"Procedural Graph":
"""

```

Figure 10: The adopted prompt consisting of elaborate instruction and three examples.