# EIPE-text: Evaluation-Guided Iterative Plan Extraction for Long-Form Narrative Text Generation

**Wang You**[*][†]   **Wenshan Wu**[*][‡]   **Yaobo Liang**[*]   **Shaoguang Mao**
**Chenfei Wu**   **Maosong Cao**[†]   **Yuzhe Cai**[†]   **Yiduo Guo**[†]   **Yan Xia**   **Furu Wei**   **Nan Duan**

Microsoft Research Asia

## Abstract

Plan-and-Write is a common hierarchical approach in long-form narrative text generation, which first creates a plan to guide the narrative writing. Following this approach, several studies rely on simply prompting large language models for planning, which often yields suboptimal results. In this paper, we propose a new framework called Evaluation-guided Iterative Plan Extraction for long-form narrative text generation (EIPE-text), which extracts plans from the corpus of narratives and utilizes the extracted plans to construct a better planner. EIPE-text has three stages: plan extraction, learning, and inference. In the plan extraction stage, it iteratively extracts and improves plans from the narrative corpus and constructs a plan corpus. We propose a question answer (QA) based evaluation mechanism to automatically evaluate the plans and generate detailed plan refinement instructions to guide the iterative improvement. In the learning stage, we build a better planner by fine-tuning with the plan corpus or in-context learning with examples in the plan corpus. Finally, we leverage a hierarchical approach to generate long-form narratives. We evaluate the effectiveness of EIPE-text in the domains of novels and storytelling. Both GPT-4-based evaluations and human evaluations demonstrate that our method can generate more coherent and relevant long-form narratives. Our code will be released in the future.

## 1 Introduction

Large language models have made impressive strides in text generation, performing well in tasks such as machine translation, summarization, and chat (Chang et al., 2023)(Bubeck et al., 2023). However, generating long-form narrative remains a challenging task, especially when it comes to maintaining coherence over long ranges and ensuring relevance to an initial premise. This is particularly crucial for applications such as scriptwriting, novels, business reports, journalism, among others.

Human writers often create a plan or outline before beginning to write a narrative, which helps maintain a coherent and logical progression throughout the narrative. Inspired by this, a hierarchical generation approach has been used in many works, such as Re3(Yang et al., 2022), DOC(Yang et al., 2023), and recurrentGPT(Zhou et al., 2023). These works mainly focus on how to generate the full narrative based on a plan and only generate the plan by simply prompting a large language model. However, the planning ability of LLMs is not good enough and requires significant prompting engineering work. Additionally, it is challenging to adapt these models to a specific domain or style of long-form narrative.

To address these limitations, we propose the Evaluation-Guided Iterative Plan Extraction for Long-Form Narrative Text Generation (EIPE-text) framework. EIPE-text leverages a learned planner with enhanced domain expertise to generate a high-quality plan, as illustrated in figure 1. Specifically, EIPE-text consists of three stages: plan extraction, learning, and inference. In the plan extraction stage, we iteratively extract and improve plans from collected narrative corpus to construct a plan corpus for planner learning. To evaluate the quality of extracted plans and the alignment between plans and source narratives, we adopt a QA-based self-evaluation mechanism, leveraging the reading comprehension capabilities of LLMs. Based on evaluation results, we generate detailed refinement instructions to iteratively improve the plan. In the learning stage, we build a better planner by fine-tuning with the plan corpus or in-context learning with examples in the plan corpus, which enhances the ability to generate high-quality plans. During

---

[*]Equal contribution
[†]Work done during internship at Microsoft Research Asia.
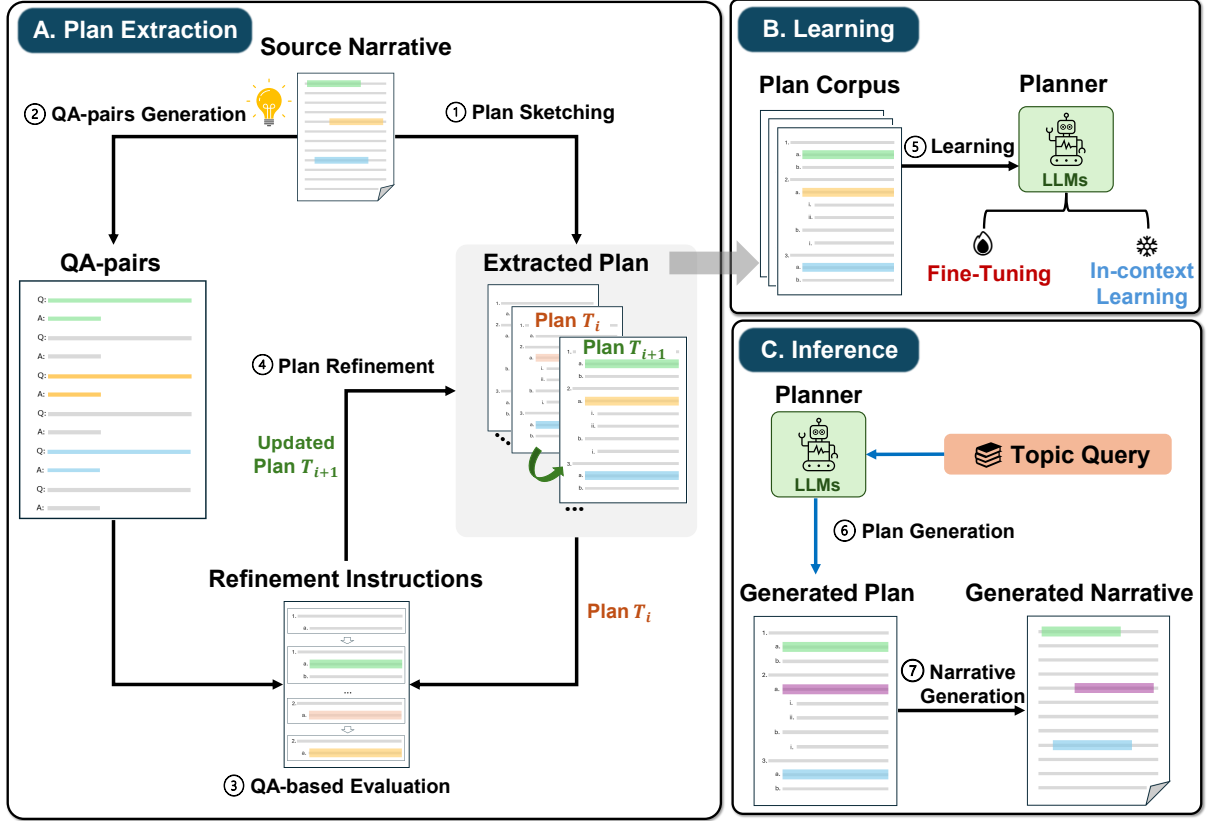[‡]Corresponding author: wenswu@microsoft.com

Figure 1: **A Comprehensive Visual Overview of the EIPE-text Framework.** The **Plan Extraction** stage starts with *Plan Sketching*, where an initial plan is generated using an LLM. Then, in the *QA-pairs Generation* step, a set of QA-pairs is created to evaluate the plan. *QA-based Evaluation* step evaluates the plan through question answering and generates refinement instructions. In the *Plan Refinement* step, it iteratively improves the plan based on the instructions until it passes the evaluation. Plans are then used to construct a plan corpus for the planner in the **Learning** stage. Finally, in the **Inference** stage, the planner generates a plan, and the narrative is generated from the plan.

the inference stage, we first generate the plan and then further generate narratives based on the plan.

We evaluated the effectiveness of EIPE-text in the domain of novels and storytelling and found that both the fine-tuning based and in-context learning based planners outperform the baselines. Human evaluation also shows that the results of EIPE-text were more coherent and relevant than those of current state-of-the-art models.

Our contributions can be summarized as follows:

- We propose a new framework, EIPE-text, which automatically extracts high-quality plans from narrative corpus and learns better planners for long-form narrative text generation. This framework can be generalized to all domains.

- We propose a QA-based evaluation method to automatically evaluate plans and generate detailed instructions to improve the plan based

on evaluation results. This QA-based evaluation provides more specific and actionable results than simply leveraging GPT to compare two outputs or provide a score (Liu et al., 2023).

- We demonstrate the effectiveness of our model in the novel and storytelling domains, and we will release the code for future research.

## 2 Method

Our methodology contains three stages: plan extraction, learning, and inference. The entire process is shown in figure 1. During the plan extraction phase, plans are extracted from each narrative within the corpus. These extracted plans are then compiled to construct the plan corpus. By relying on the constructed planning corpus, planner can learn to generate high-quality plans. In the inference stage planner generates a better plan and a

narrative will be generated from the plan.

The plan extraction stage contains plan sketching, QA-pairs generation, QA-based evaluation, and plan refinement. Initially, we create a tree-structured plan using the LLM in the plan sketching step. Next, during the QA-pairs generation phase, we generate a set of QA-pairs, with each pair corresponding to a distinct part within the source narrative. These QA-pairs serve as an evaluation metric for the plan. The QA-based evaluation step evaluates the plan by question answering. For each incorrect QA-pair, we generate corresponding instructions to modify the relevant part of the plan. In the plan refinement step, we integrate these instructions received in previous steps to update the plan. We repeat steps 3 and 4 until the extracted plan passes the evaluation.

In the learning stage, We leverage the plan extracted in the first stage to train an LLM planner. To achieve this, we utilize two strategies: finetuning, as well as in-context learning. These strategies contribute to generating high-quality plans for the given topic.

The inference stage contains two steps: plan generation and narrative generation. Firstly, the planner takes the topic as input and generates a corresponding plan. Secondly, the narrative will be generated in the narrative generation step.

## 2.1 Plan Extraction

Formally, we have a corpus of narrative $\mathcal{C}_n = \{n_1, n_2, ..., n_m\}$. The plan extraction stage extracts a plan $p_i$ for each narrative $n_i$. The extraction results are compiled to a plan corpus $\mathcal{C}_p = \{p_1, p_2, ..., p_m\}$. We illustrated the process of plan extraction in algorithm 1.

**Plan Sketching.** For each narrative, we use LLM to extract a tree-structured plan, which serves as the plan sketch. The detailed LLM prompt can be found in appendix A.1. The plan is in a tree structure and the content of each **node** is the summarization of the corresponding section, subsection, and so forth. We show an example of a plan sketch in figure 2.

**QA-pairs Generation.** For each narrative, we generate a set of QA-pairs, with each pair corresponding to a different segment of the narrative. These QA-pairs can be utilized to evaluate whether the plan includes all aspects of the narrative. Each QA-pair is formulated as a multiple-choice problem, comprising one question, multiple candidate

answers, and multiple correct answer indices. The number of QA-pairs is proportional to the length of the narrative. To ensure the quality of the generated QA-pairs, we employ another LLM to answer these questions based on the original text, filtering out any incorrectly answered pairs. The guidelines for this process can be found in appendix A.2.

**QA-base Evaluation.** We evaluate a plan using QA-pairs and provide detailed refinement instructions for refining the plan further. Specifically, we utilize LLM to answer questions based on the plan. For each incorrect question, we generate an instruction to modify the plan so that the question can be correctly answered. The modification instruction can be one of the following: (1) **add**, which inserts a missing node into the plan; (2) **modify**, which alters the content of a node; (3) **adjust**, which relocates a node to another level of the tree, thereby altering the tree's structure. Detailed refinement instructions enable LLM to make precise improvements to specific parts of the plan.

**Plan Refinement.** In this step, we incorporate the instructions generated in the previous step to improve the plan. Ideally, we should apply the changes one by one. In order to improve efficiency, we instruct the LLM to apply all instructions simultaneously. However, the refinement instructions generated by LLM may not always address the incorrect questions. Therefore, we iteratively perform the refinement instructions generation and plan refinement steps until the new plan can pass the QA-based evaluation. This process ensures that the final plan has addressed all the identified errors and meets the desired quality standards.

While LLM possesses a self-improving ability and can refine the plan through simple prompting, the quality of the improvement results may still not be good enough or even worse. Our QA-based evaluation, on the other hand, can identify specific errors in the plan and provide refinement instructions in the form of instructions to enhance the plan. This approach can achieve better refinement performance.

## 2.2 Learning

During the learning phase, we implemented two methods to enhance the performance of the planner: the in-context learning method and the fine-tuning method.

The in-context learning method improves the planner by selecting representative demonstration

---

**Algorithm 1:** Plan Extraction Algorithm

---

**Input** : $\mathcal{C}_n = \{n_1, n_2, ..., n_m\}$

**Output** : $\mathcal{C}_p = \{p_1, p_2, ..., p_m\}$

1   $\mathcal{C}_p \leftarrow \varnothing$

2   **for** $i \leftarrow 1$ **to** $m$ **do**

3      $p_i^0 \leftarrow$ plan_sketching$(n_i)$

4      $\mathcal{C}_q \leftarrow$ qa_pairs_generation$(n_i)$              $\triangleright \mathcal{C}_q = \{q_1, q_2, ..., q_k\}$ **questions set**

5      $t \leftarrow 0$                                           $\triangleright t$ **refinement time step**

6      **while** *not pass_evaluation*$(p_i^t, \mathcal{C}_q)$ **do**

7         $\mathcal{C}_i \leftarrow$ qa_based_evaluation$(p_i^t, n_i, \mathcal{C}_q)$    $\triangleright \mathcal{C}_i = \{i_1, i_2, ..., i_l\}$ **refinement instructions set**

8         $p_i^{t+1} \leftarrow$ plan_refinement$(p_i^t, \mathcal{C}_i)$

9         $t \leftarrow t + 1$

10      **end**

11      $\mathcal{C}_p \leftarrow \mathcal{C}_p \cup p_i^t$

12 **end**

---

examples from the plan corpus. By selecting different demonstration examples, the fixed LLM can quickly adapt to specific domains or styles.

On the other hand, the fine-tuning method can further improve the planner's ability by training it on all plan corpus. This method leverages all the data in the plan corpus and enables the planner to adapt to multiple domains simultaneously.

## 2.3 Inference

The inference stage comprises two steps: plan generation and narrative generation.

**Plan Generation.** In this step, the planner takes the chosen topic as input and produces a corresponding plan. The planner constructs a well-structured plan that outlines the key elements and sections to be covered in the ensuing narrative.

**Narrative Generation.** The narrative is generated from the generated plan in this step. This narrative seamlessly integrates the content outlined in the plan, ensuring that the resulting narrative is not only logically organized but also rich in detail and context. The final narrative is a well-rounded piece of long-form narrative that effectively conveys the information related to the chosen topic.

## 2.4 Discussion

In this section, we will discuss how EIPE-text works. Here is our analysis:

Let $q$ be the premise query. The probability of desired output based on premise query $p(n|q)$ could be rewritten as

$$P(n|q) = P(p|q)P(n|p) \tag{1}$$

When plan $p$ is of high quality, $P(n|p)$ will be high. So as $P(p|q)$ increases, $P(n|q)$ increases too. Our framework EIPE-text actually increases $P(p|q)$.

Besides, the process of plan refinement in figure 1 could be understood as Reinforcement Learning(RL), LLM gets observation from answering the question, and then obtains refinement instructions according to the true or false case. After obtaining refinement instructions, LLM changes the original state to the new state i.e. revise plan. After many interactions with the "environment", the "state" will be iterated to a suitable "state" that can be used to improve $P(p|q)$.

To practically exemplify the effectiveness of EIPE-text, we conducted a case study of plan generation through in-context learning with one demonstration. A detailed exploration of this case is provided in the Appendix D.2 for interested readers.

## 3 Experiments

In this section, we compare EIPE-text in novels and storytelling generation with the baselines. All experiments show that EIPE-text is better than the baselines, verifying the effectiveness of our framework.

## 3.1 Setup

For plan extraction stage, we use Azure Openai GPT-4 as our experimental LLM. And for inference stage, we use the planner to generate a plan to further generate the narrative. It should be emphasized that we did not intentionally implement the narrative generation, but modified it based on recurrentGPT, as described in the appendix B.1.

| Dataset | Train Size | Test Size | Avg Length | Max Length |
|---------|-----------|-----------|-----------|-----------|
| TED Talk | 2468 | 130 | 2078 | 9044 |
| Novel | 1292 | 120 | 3741 | 14493 |

Table 1: Comprehensive Dataset Information for TED Talk and Novel.

| Novel genres | Overall(human) | | |
|---------|-----------|-----------|-----------|
| ~4500words | Interesting | Coherent | Relevant |
| EIPE-text (in-context) | 56.7 | **64.2** | **75.8** |
| recurrentGPT | **60.0** | 59.2 | 62.5 |

Table 2: Novel Human Evaluation Results. Pair-wise comparison using human evaluation of EIPE-text with recurrentGPT for 120 novels of different genres. Results never mix numbers from different comparisons

| Novel genres | Overall(automatic) | | |
|---------|-----------|-----------|-----------|
| ~4500words | Interesting | Coherent | Relevant |
| EIPE-text (in-context) | 55.0 | **84.2** | **92.5** |
| recurrentGPT | **58.3** | 65.8 | 84.2 |

Table 3: Novel GPT4 Evaluation Results. Pair-wise comparison using GPT-4 evaluation of EIPE-text with recurrentGPT for 120 novels of different genres. Results never mix numbers from different comparisons

**For all the settings mentioned in the following section, unless special emphasis, they adhere to the description provided above.**

### 3.2 Novel

#### 3.2.1 Dataset

Novels are long-form narratives that include intricate plots, and rich character development. The model needs to maintain consistency in plots and character development and generate interesting stories. We use the data collected from Project American Literature[1], Writing Prompts[2] and etc. Then we aggregate a training dataset containing total 1292 stories. Besides, we collected 120 prompts as a test set from Writing Prompts, which cover six genres. The more information about this dataset is shown in table 1.

#### 3.2.2 Setting

**EIPE-text (in-context)** For learning stage, we use the *text-embedding-ada-002*, to obtain text embeddings of plan corpus. These embeddings will

---

[1] https://americanliterature.com/short-stories
[2] https://blog.reedsy.com/creative-writing-prompts/

then be utilized in conjunction with the *k-means* algorithm for cluster purposes. We use *k-means* getting 20 clustering centroids as demonstrations to learn a planner and use the planner during comparing with baselines.

#### 3.2.3 Baselines

**recurrentGPT** A language-based simulacra of the recurrence mechanism in RNNs that uses language-based components and defines a recurrent computation graph via prompt engineering.

It is worth mentioning that we are not directly comparing with Re3 and DOC, because recurrentGPT is already way ahead of these methods.

#### 3.2.4 Metric

Our evaluation employs a pairwise comparison metric. We report results individually for each pairwise comparison between EIPE-text and each baseline, never mixing numbers from different comparisons following Re3 (Yang et al., 2022). We show the criteria as outlined in (Yang et al., 2023) for novel as following:

- **Interesting**: An interesting novel captivates the reader's attention, engages them emotionally, and holds their interest throughout.

- **Coherent**: A coherent novel follows a logical and consistent plot-line without significant gaps or inconsistencies.

- **Relevant**: Faithful to the initial premise.

**Automatic Evaluation** For automatic evaluation, we employed GPT-4 to assess various aspects of the generated narrative. GPT-4 automatic evaluation is highly affected by the order and unstable, so all metrics are judged by GPT4 with a premise, aforementioned criteria and two corresponding stories in random order. We also use majority voting system to evaluate each criterion of each pair. The evaluation prompt for novel can be found in appendix C.1.

| setting A | setting B | A Win Ratio | B Win Ratio |
|---|---|---|---|
| LLaMA raw planner | EIPE-text (finetune) | 6.2 | **93.8** |
| GPT4 raw planner | EIPE-text (in-context) | 22.5 | **75.2** |

Table 4: TED Talk Automatic Evaluation Results. Pair-wise comparison using GPT-4 evaluation of EIPE-text with baselines for 130 TED talk transcripts. Results in different comparisons are not comparable with each other.

**Human Evaluation**　In order to ensure impartial and high-quality evaluations, we collaborated with third-party data annotators. Each generated data pair, comprising novels A and B presented in random order, underwent meticulous evaluation by three distinct annotators. These annotators possess proficient English language skills and were provided with explicit instructions to evaluate and deliver judgments on the superiority between novel A and novel B, or if they are indistinguishable, specifically in relation to the aforementioned criteria.

### 3.2.5　Result

We show the experiment results of novels in table 2 and table 3. As we can see from the table, EIPE-text shows an advantage in coherence and relevance in both human and automatic evaluation. Although the human evaluation is less interesting (3.3%), the improvement of coherence (5.0%) and relevance (13.3%) are significant. The same trend can be seen in automatic evaluation, it is less interesting than recurrentGPT(3.3%), but coherent (18.4%) and relevant (8.3%) are significantly higher. These results indicate that EIPE-text improves the overall quality of generated narrative, and also indicate that automatic evaluation and human evaluation have certain relevance.

### 3.3　Strorytelling

### 3.3.1　Dataset

TED Talks [3] are influential presentations that cover a wide range of topics. They are known for their engaging narratives, concise structure, and powerful messages, which can be challenging to generate for both models and humans. We use the data collected by Kaggle [4]. The training dataset aggregates 2,468 TED Talks spanning the years 1984 to 2016. In addition, we have curated 130 TED Talk transcripts post-2021 as our testing datasets as shown in table 1.

---

[3] https://www.ted.com/talks
[4] https://www.kaggle.com/datasets/rounakbanik/ted-talks

### 3.3.2　Setting

**EIPE-text (in-context)**　For learning stage, text embeddings obtained using *text-embeddings-ada-002* are used for clustering together with the *k-means* algorithm. Then we use 20 clustering centroids as demonstrations to learn a planner.

**EIPE-text (finetune)**　We finetune the open source LLM, LLaMA (Touvron et al., 2023), using the plan corpus and use it as planner during learning stage. Specially, we finetune LLaMA-7B using lora(Hu et al., 2022).

### 3.3.3　Baselines

**GPT4 raw planner**　In this setup, planner is GPT4 zero-shot whose ability to plan depends entirely on its native capabilities. After the planner generates the plan, narrative generation follows the same way as the inference stage in 3.1

**LLaMA raw planner**　similar to GPT4 raw planner, but the planner is untrained LLaMA.

### 3.3.4　Metric

We only adopt automatic evaluation in storytelling generation. The evaluation criteria were tailored to specific domain to ensure relevant and accurate assessments, so we use other criteria for storytelling:

- **Coherent**: The talk should have a clear structure, smooth transitions, and a strong conclusion for easy comprehension and a consistent theme.

- **Interesting**: It should use storytelling and examples to engage the audience, maintaining their curiosity throughout.

- **Relevant**: The topic should be timely, address current issues, and offer fresh insights, not just repeat existing information.

- **Inspiring**: The talk should convey passion, present innovative ideas, and encourage the audience to think differently or take action.

It should be emphasized that we only use majority voting system to evaluate each pair for all criteria, instead of evaluating each criterion of each pair. The evaluation prompt for storytelling can be found in appendix C.2

### 3.3.5 Results

We show the experiment result of storytelling domain on TED Talk in table 4. Under the finetune setting, EIPE-text far outperforms LLaMA raw planner (87.6%). Also under setting B, EIPE-text is significantly outperform the GPT4 raw planner (52.7%). EIPE-text either using a finetune base planner or using in-context learning based planners is well ahead of the LLM itself.

## 4 Analysis

In this section, we explore the key aspects of designing an effective planner and provide an experimental analysis of the effectiveness of the plan refinement process.

### 4.1 Ablation study of in-context learning based planner

Our investigation centers around two fundamental questions: (1) How does the demonstration selection algorithm impact the performance of our planner? (2) What effect does the number of demonstration examples have on the planner's performance?

To address these questions, we designed experiments where we compared various planner configurations, including (1) **n-shot cluster-based planner**: this configuration utilizes a cluster-based approach to select n demonstration examples. (2) **n-shot retrieval-based planner**: in contrast, this configuration employs a retrieval-based method to select n demonstration examples.

**Using clustering to select more demonstrations leads to better results.** We show the results in table 5. In the comparison between the 20-shot cluster-based planner and the 5-shot cluster-based planner, the 20-shot cluster-based planner outperforms the 5-shot cluster-based planner with a win ratio of 70.9% versus 26.8%. This suggests that using more demonstration examples leads to better planner performance. In addition, as the plan length we use is shorter than full narrative, we can use more plans as demonstrations within context window. When comparing the 5-shot cluster-based planner and the 5-shot retrieval-based planner, the clustering-based method for selecting demonstration examples appears to be slightly more effective.

This trend is more pronounced when looking at the comparison between the 20-shot cluster-based planner and the 20-shot retrieval-based planner. The 20-shot cluster-based planner significantly outperforms the retrieval-based planner, with a win ratio of 67.2% versus 32.0%. This suggests that using clustering for selection is considerably more effective than relying on retrieval-based methods.

### 4.2 Comparison between hierarchical generation and non-hierarchical

To investigate the impact of narrative generation methods on the performance of our planner, we compared hierarchical generation with non-hierarchical methods.

We experiment with non-hierarchical generation including configurations: (1) **0-shot without planner**: generate full narrative directly in one step. (2) **n-shot cluster-based without planner**: select n demonstrations using a cluster-based approach and generate a full narrative using these demonstrations. (3) **n-shot retrieval-based without planner**: similar to previous setting, instead, we rely on a retrieval-based approach to select demonstrations.

**Hierarchical generation is effective compared with non-hierarchical**. We show the results in table 5. The 0-shot planner, significantly outperforms 0-shot without planner, achieving a win ratio of 76.7% versus 20.9%. Moreover, similar trends can be found in 5-shot setting with 88.2% versus 11.0% and 70.6% versus 29.4%.

### 4.3 Effectiveness of the plan refinement process

In addition, we also want to know whether self-refinement can be effectively refined and the reasons behind its convergence.

**Fast Convergence with Self-Refinement** We can see from the table 6 that our framework can converge in an average of 2.98 epochs, which is actually very fast and it is hard to converge without using self-refinement. The average accuracy curve of iterative refinement process is shown in figure 3.

**Iterative Plan Refinement Ensures Alignment** The refined plan contains three operations, we monitor the number of three operations in the process. In addition, since we organize the plan into a tree structure, we also record the change in the number of nodes in the tree and the change in the number of secondary nodes (children of the root node) throughout the process. As can be seen from table

Figure 2: An Example of the Plan Refinement Process.



Figure 3: Average accuracy curve of iterative refinement process.

6, the average add, modify and adjust operations occur 8.26 times, 3.22 times, and 2.25 times respectively. The average number of nodes increase by 11.41. We can clearly see these changes in figure 2 (for more detail in appendix D.1). This indicates that in plan refinement process, it does not simply add nodes. Instead, it can accurately modify relevant parts and adjust structure according to the question answering. Thus, these three operations ensure the alignment between the plan and the original narrative.

## 4.4 Case study of in-context learning based plan generation

Relying solely on comprehensive narratives for learning can often lead to missing finer details. Narratives are typically dense with information, posing challenges for models to pinpoint and retain crit-

| A | B | A Win Ratio | B Win Ratio |
|---|---|---|---|
| **Different Demonstration Number** | | | |
| 20-shot cluster-based planner | 5-shot cluster-based planner | **70.9** | 26.8 |
| **Different Demonstration Selection** | | | |
| 5-shot cluster-based planner | 5-shot retrieval-based planner | **51.6** | 46.0 |
| 20-shot cluster-based planner | 20-shot retrieval-based planner | **67.2** | 32.0 |
| **Different Narrative Generation Method** | | | |
| 0-shot planner | 0-shot without planner | **76.7** | 20.9 |
| 5-shot cluster-based planner | 5-shot cluster-based without planner | **88.2** | 11.0 |
| 5-shot retrieval-based planner | 5-shot retrieval-based without planner | **70.6** | 29.4 |

Table 5: **Ablation Study Result. Different Demonstration Number**: In the learning stage of EIPE-text, in-context learning based planner use different numbers of demonstrations. **Different Demonstration selection**: In-context learning based planner can implement different methods, such as clustering or retrieving items related to the input topic, to select demonstrations. **Different Narrative Generation Method**: In addition to being able to generate narratives using EIPE-text. Narrative can also be generated in one step by simply combining several narratives as demonstrations without planner giving an input topic.

| metric | operation | | | difference before and after | | epochs and question numbers | |
|---|---|---|---|---|---|---|---|
| | **add** | **modify** | **adjust** | **all nodes** | **secondary nodes** | **average epoch** | **average questions** |
| num | 8.26 | 3.22 | 2.25 | 11.41 | 0.25 | 2.98 | 35.71 |

Table 6: Iterative Refinement Metric

ical elements. Furthermore, methods that learn from complete narratives are usually computationally expensive and time demanding. On the other hand, when using in-context learning with plans, models can more adeptly identify and relate to relevant information within each contextual segment. This technique not only ensures that key details aren't overlooked but also streamlines the learning process regarding the text's semantic framework, ultimately conserving computational resources. We show an example of 1-shot in Appendix D.2, from which we can see that the generated plan is not only coherent but also retains the salient features of the demonstration, while effectively addressing the topic query.

## 5 Related Work

**Long-form Narrative Text Generation** As for long-form narrative text generation, recent studies tackle this from the following perspectives: appending the generated prefix to the encoder (Shao et al., 2017), while newer models like (Guan et al., 2021) focus on capturing sentence and discourse-level coherence, and DiscoDVT by (Ji and Huang, 2021) leverages discrete variational Transformers to enhance long-range coherence in generated texts. Another type of work adopts the plan-and-write strategy (Fan et al., 2018). In particular, there has been extensive exploration of story planning (Yao

et al., 2019; Fan et al., 2019; Goldfarb-Tarrant et al., 2020). A hierarchical story generation system with recursive prompting and revision was proposed by Yang et al. (2022). And the current state-of-the-art work recurrentGPT (Zhou et al., 2023), which uses large language model (LLM) such as ChatGPT and uses natural language to simulate the Long Short-Term Memory mechanism in an LSTM. The current plan results from these methods are not satisfactory. Instead, we use LLM to automatically mine the plan and train a good planner to achieve good results. Furthermore, from the plan to the full text, our methods and theirs are complementary and can be combined to achieve better results.

**Human-AI Co-writing** Human-AI co-writing systems have been developing at the intersection of NLP and human-computer interaction (HCI) fields, such as Wordcraft (Yuan et al., 2022), TaleBrush (Chung et al., 2022), CoAuthor (Lee et al., 2022) and Dramatron (Mirowski et al., 2023). These works explore the possibilities of using LLM as a writing assistant to humans. Our work generates an explicit plan, which can be easily provided for human review and modification, making human-AI co-writing easier.

# 6    Conclusions

EIPE-text represents a significant step forward in the field of long-form narrative text generation, addressing the challenges of coherence and structure over extended pieces of text. With its ability to generate high-quality long-form narratives and aid human writers, EIPE-text opens up new possibilities for leveraging the capabilities of LLMs in creative and expressive writing tasks. Future research could explore further applications and extensions of EIPE-text in various domains, advancing the state of the art in automated text generation.

# 7    Limitations

During plan extraction stage, the two steps of QA-pairs generation and questions answering largely depend on LLM's own reasoning capability, so this method can only produce ideal results on models with strong reasoning capability (GPT4, Claude, etc.). Otherwise, it may lead to the refinement process failing to converge. Our framework is a data-driven approach, so it does not improve the OOD performance.

# References

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2023. A survey on evaluation of large language models.

John Joon Young Chung, Wooseok Kim, Kang Min Yoo, Hwaran Lee, Eytan Adar, and Minsuk Chang. 2022. Talebrush: sketching stories with generative pretrained language models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–19.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.

Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. In *Proceedings of the 57th Annual Meeting of the Asso-ciation for Computational Linguistics*, pages 2650–2660, Florence, Italy. Association for Computational Linguistics.

Seraphina Goldfarb-Tarrant, Tuhin Chakrabarty, Ralph Weischedel, and Nanyun Peng. 2020. Content planning for neural story generation with aristotelian rescoring. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4319–4338, Online. Association for Computational Linguistics.

Jian Guan, Xiaoxi Mao, Changjie Fan, Zitao Liu, Wenbiao Ding, and Minlie Huang. 2021. Long text generation by modeling sentence-level and discourse-level coherence. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6379–6393, Online. Association for Computational Linguistics.

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Haozhe Ji and Minlie Huang. 2021. DiscoDVT: Generating long text with discourse-aware discrete variational transformer. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4208–4224, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Mina Lee, Percy Liang, and Qian Yang. 2022. Coauthor: Designing a human-ai collaborative writing dataset for exploring language model capabilities. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–19.

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. Gpteval: Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*.

Piotr Mirowski, Kory W Mathewson, Jaylen Pittman, and Richard Evans. 2023. Co-writing screenplays and theatre scripts with language models: Evaluation by industry professionals. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–34.

Yuanlong Shao, Stephan Gouws, Denny Britz, Anna Goldie, Brian Strope, and Ray Kurzweil. 2017. Generating high-quality and informative conversation responses with sequence-to-sequence models. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2210–2219, Copenhagen, Denmark. Association for Computational Linguistics.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal

Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.

Kevin Yang, Dan Klein, Nanyun Peng, and Yuandong Tian. 2023. DOC: Improving long story coherence with detailed outline control. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3378–3465, Toronto, Canada. Association for Computational Linguistics.

Kevin Yang, Yuandong Tian, Nanyun Peng, and Dan Klein. 2022. Re3: Generating longer stories with recursive reprompting and revision. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4393–4479, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7378–7385.

Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. 2022. Wordcraft: story writing with large language models. In *27th International Conference on Intelligent User Interfaces*, pages 841–852.

Wangchunshu Zhou, Yuchen Eleanor Jiang, Peng Cui, Tiannan Wang, Zhenxin Xiao, Yifan Hou, Ryan Cotterell, and Mrinmaya Sachan. 2023. Recurrentgpt: Interactive generation of (arbitrarily) long text.

## A  Prompts

### A.1  Prompt of Sketching Plan

Distill the salient information and thematic flow from the original article into a tree-like text representation of a mind map in the following format:

```
TOPIC
    - Main Topic
        - Sub Topic
            - Sub-Sub Topic
            - Sub-Sub Topic
    ...
    - Main Topic
        - Sub Topic
        - Sub Topic
```

### A.2  Prompt of QA-pairs Generation Guideline

Based on the content of the article, generate several multiple-choice questions and corresponding answers:

1. Not too detailed

2. Focus on the logic of the article

3. Deep understanding of the article after answering these questions

4. Each question must have 4 options: A, B, C, D.

5. For each question, there might be more than one correct answer, identify all correct answers separated by ";"

6. Questions should reflect the structure of the article.

7. Questions should include three types: what, why, how.

8. Provide related main ideas in the article for each question.

9. Avoid options like "All of the above" or "None of the above"; use "A;B;C" format.

These questions are generated based on the article's content and the author's opinion, not my opinion.

## B  Experiment Details

### B.1  Modification for recurrentGPT

The way to improve recurrentGPT. recurrentGPT is prone to loss of global memory just as RNN. And we also find that the long-term memory in recurrentGPT is not exactly long-term memory. To compensate for this, we can insert the generated plan as additional memory to recurrentGPT, as shown in figure 4. At each time step ('t'), recurrentGPT operates on a dual input system: the paragraph produced in the preceding step and a concise yet directive instruction for the subsequent paragraph. A crucial aspect is the integration of the model's long-term memory, which acts as a repository for storing previously generated summaries, and importantly, it can retrieve these summaries through semantic search, with the ability to store them on external hard drives. Simultaneously, the system actively maintains a short-term memory, responsible for encapsulating key information from recent time steps, a repository that gets updated consistently as the process unfolds. Crucially, the "generated plan," a newly introduced memory, becomes an integral part of this intricate orchestration. When the components converge, they create a coherent prompt that triggers the backbone language model, aptly dubbed the "backbone LLM," to undertake its
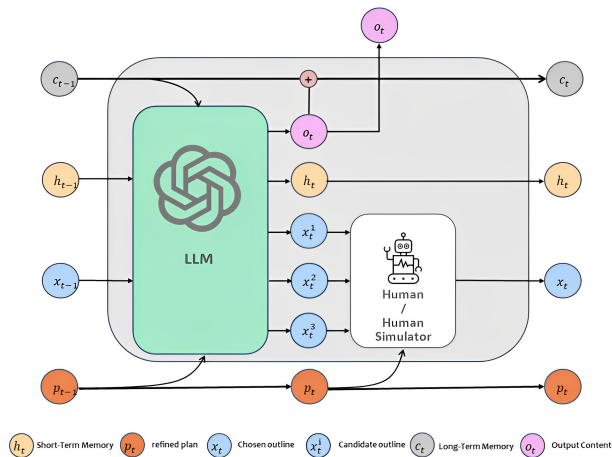
Figure 4: Illustration of the improvement of recurrent-GPT. recurrentGPT uses natural language simulating an LSTM, we insert additional memory to main the long-term memory in LSTM

primary task: generating a fresh paragraph while simultaneously outlining a succinct plan for the forthcoming paragraph. What's truly remarkable is how this "generated plan" seamlessly merges with the process, as it is not only updated in each time step but also contributes to enriching the long-term memory. This meticulous integration ensures continuity and coherence throughout the sequence, forging a recurrent mechanism that drives the generation process forward, where the "generated plan" plays a pivotal role in shaping the narrative's development.

### B.2 Use K-Means Get Demonstrations

Using text-embeddings-ada-002 directly to convert a text format plan to embedding for clustering is not varied. We have adopted a combination approach of K-means and LLM. To be specific, we use LLM, according to prompt "Without loss of generality, list distinctive characteristics of this exemplar that establish it as an effective paradigm for designing genre. no explanation is needed." to get characteristics of plans. Then these characteristics are converted into embedding before clustering. After we set the number of clusters k, we can get k clustering centers, and we use these centroid as final demonstrations.

## C Automatic Evaluation

### C.1 Novel Automatic Evaluation

In this task, you will be presented with two novels side-by-side and asked to evaluate them based on three metrics: Coherence, Interestingness, Rele-

vance. Your task is to determine which novel is better for each metric or indicate if both novels are indistinguishable.

- Coherent: A coherent novel follows a logical and consistent plot-line without significant gaps or inconsistencies.

- Interesting: An interesting novel captivates the reader's attention, engages them emotionally, and holds their interest throughout.

- Relevant. Faithful to the initial premise. The novel effectively aligns its plot, message, and writing with its initial premise, ensuring consistency and faithfulness to the core theme.

Based on these three aspects, make a decision on which novel is achieving the desired impact, in the manner like this:

[Scratch Pad]

Name:

'distinctive characteristics' and 'elaborate' on them

Name:

'distinctive characteristics' and 'elaborate' on them

[Reflection]

After evaluating both novels based on the criteria of 'coherence', 'interestingness', 'relevance', I have come to the following 'thorough' conclusions:

- Coherence:

- Interestingness:

- Relevance:

[Final Choice]:

Coherence: Name;

Interestingness: Name;

Relevance: Name;

### C.2 TED Talks Automatic Evaluation

The coach's preference for evaluating the TED Talks can be summarized in the following spec:

- Coherence: The coach will assess how well the TED Talk is structured and organized. This includes a clear introduction, logical flow of ideas, smooth transitions between points, and a strong conclusion. The talk should be easy to follow and understand, with a consistent theme throughout.

- Interestingness: The coach will evaluate how engaging and captivating the TED Talk is for the audience. This includes the use of storytelling, anecdotes, and examples to illustrate points, as well as the speaker's ability to maintain the audience's attention and curiosity throughout the talk.

- Relevance: The coach will consider the importance and significance of the topic being discussed

in the TED Talk. The subject matter should be timely, relevant to current events or societal issues, and have a broad appeal to a diverse audience. The talk should also provide new insights or perspectives on the topic, rather than simply rehashing existing information.

- Inspiration: The coach will assess the TED Talk's ability to inspire, motivate, and provoke thought in the audience. This includes the speaker's ability to convey passion and enthusiasm for the topic, as well as the presentation of innovative ideas, solutions, or calls to action that encourage the audience to think differently or take action in their own lives.

Based on these four aspects, the coach will make a decision on which TED Talk is stronger and more effective in achieving the desired impact on the audience, in the manner like this:

[Scratch Pad]

Name:

'distinctive characteristics' and 'elaborate' on them

Name:

'distinctive characteristics' and 'elaborate' on them

[Reflection]

After evaluating both TED Talks based on the criteria of 'coherence', 'interestingness', 'relevance', and 'inspiration', I have come to the following 'thorough' conclusions:

- Coherence:
- Interestingness:
- Relevance:
- Inspiration:

[Final Choice]: Name

## D Examples

### D.1 Plan Extraction Example

In this section, we show the detailed process of an iteration in the plan extraction stage in the figure 5. And the comparison of initialized plans and refined plans are also shown in figure 6, 7, 8

### D.2 Plan Generation Example

In figure 9, we observe several salient aspects of the demonstration on the left side: 1. the utilization of relatable examples and analogies. 2. the connection of various disciplines and concepts. 3. the incorporation of quotes from notable figures. These attributes are integrated into the generated plan. Notably, the generated plan contains relatable examples to substantiate its viewpoints and incorporates relevant quotes. When we transition to the generated plan on the right, it's evident that the plan incorporates these attributes seamlessly. For instance: 1. the mention of "Daniel Kahneman's System 1 and System 2 thinking" in the plan mirrors the demonstration's theme on "Human perception and understanding.". 2. the outcome's emphasis on quotes is reflective of the demonstration's approach, incorporating wisdom from Peter Drucker and Tim Ferriss. We can observe that the generated plan not only maintains coherence but also preserves the key features of the demonstration, while effectively responding to the topic query.

**Initialized Plan**

Pig Products in Daily Life
  - Introduction
    - Netherlands: 16 million people, 12 million pigs
    - Research on pig usage in products
  - Products containing pig parts
    - Bathroom items
      - Soap (fatty acids from pork bone fat)
      - Shampoo, conditioner, anti-wrinkle cream, body lotion, toothpaste
    - Food items
      - Dough improver (proteins from pig hairs)
      - Low-fat butter (gelatin for texture)
      - Cheesecake, chocolate mousse, tiramisu, vanilla pudding (gelatin for appearance)
    - Construction materials
      - Cellular concrete (proteins from bones)
      - Train brakes (bone ash)
    - Household items
      - Fine bone china (translucency and strength)
      - Paint (texture and glossiness)
      - Sandpaper (bone glue)
      - Paintbrushes (pig hairs)
    - Meat products
      - Portion-controlled meat cuts (fibrin from pig blood)
    - Beverages
      - Beer, wine, fruit juice (gelatin for clarity)
    - Other products
      - Cigarettes with hemoglobin filters
      - Injectable collagen for wrinkles
      - Bullets
      - Heart valve implants
      - Renewable energy (fuel from unused pig parts)
  - Conclusion
    - 185 products found containing pig parts
    - Importance of knowing what products are made of
    - Taking better care of raw materials and producers

**QA-based evaluation**

**Question 1:** Why does soap contain fatty acids made from boiling pork bone fat?

- A. To harden the soap
- B. To give it a pearl-like effect
- C. To make it smell better
- D. To improve its cleaning ability

**Answer 1:** A; B

**Wrong Answer 1:** A

**Reference 1:** "Level 4 node 'Soap (fatty acids from pork bone fat)'"

**Question 2:** What is the purpose of using pig hemoglobin in cigarette filters?

- A. To create an artificial lung
- B. To improve the taste of the cigarette
- C. To make the cigarette burn slower
- D. To reduce the amount of harmful chemicals

**Answer 2:** A

**Wrong Answer 2:** D

**Reference 2:** "Level 4 node 'Cigarettes with hemoglobin filters'"

**Question 3:** What is the purpose of using pig proteins in cellular concrete?

- A. To strengthen the concrete
- B. To make the concrete lighter
- C. To improve the concrete's heat resistance
- D. To increase the concrete's elasticity

**Answer 3:** B

**Wrong Answer 3:** A

**Reference 3:** "Level 4 node 'Cellular concrete (proteins from bones)'"

**Question 4:** Why did the director of the heart valve company not want their product to be associated with pigs?

- A. Because pigs are considered unclean animals
- B. Because it would reduce the perceived value of the product
- C. Because they were concerned about potential religious objections
- D. Because they wanted to maintain the high-tech image of the product

**Answer 4:** D

**Wrong Answer 4:** C

**Reference 4:** "Level 4 node 'Heart valve implants'"

**Question:** How are pig-derived ingredients used in renewable energy?

- A. As a substitute for fossil fuels
- B. In the production of solar panels
- C. In the production of wind turbines
- D. As a source of biofuel

**Answer:** D

**Wrong Answer:** A

**Reference:** "Level 4 node 'Renewable energy (fuel from unused pig parts)'"

**Refinement Instructions**

**Instruction**: "Add a new level 5 node 'Hardening agent and pearl-like effect' under level 4 node 'Soap (fatty acids from pork bone fat)', as the first level 5 node"

**Instruction**: "Add a new level 5 node 'Creates an artificial lung in the filter for a healthier cigarette' under level 4 node 'Cigarettes with hemoglobin filters', as the first level 5 node"

**Instruction**: "Add a new level 5 node 'Lightweight and fully reusable' under level 4 node 'Cellular concrete (proteins from bones)', as the first level 5 node"

**Instruction**: "Add a new level 5 node 'Director did not want the product associated with pigs to maintain its high-tech image' under level 4 node 'Heart valve implants', as the first level 5 node"

**Instruction**: "Change the name of level 4 node 'Renewable energy (fuel from unused pig parts)' to 'Renewable energy (fuel from pig waste)'; Add a new level 5 node 'Pig waste is converted into fuel for renewable energy' under level 4 node 'Renewable energy (fuel from pig waste)'"

**Refine plan**
Pig Products in Daily Life
  - Introduction
    - Netherlands: 16 million people, 12 million pigs
    - Research on pig usage in products
  - Products containing pig parts
    - Bathroom items
      - Soap (fatty acids from pork bone fat)
        - <mark>Hardening agent and pearl-like effect</mark>
    - Shampoo, conditioner, anti-wrinkle cream, body lotion, toothpaste
    - Food items
      - Dough improver (proteins from pig hairs)
      - Low-fat butter (gelatin for texture)
      - Cheesecake, chocolate mousse, tiramisu, vanilla pudding (gelatin for appearance)
    - Construction materials
      - Cellular concrete (proteins from bones)

           - Lightweight and fully reusable
         - Train brakes (bone ash)
     - Household items
         - Fine bone china (translucency and strength)
         - Paint (texture and glossiness)
         - Sandpaper (bone glue)
         - Paintbrushes (pig hairs)
     - Meat products
         - Portion-controlled meat cuts (fibrin from pig blood)
     - Beverages
         - Beer, wine, fruit juice (gelatin for clarity)
     - Other products
         - Cigarettes with hemoglobin filters
            - Creates an artificial lung in the filter for a healthier cigarette
         - Injectable collagen for wrinkles
         - Bullets
         - Heart valve implants
            - Director did not want the product associated with pigs to maintain its high-tech
image
         - Renewable energy (fuel from unused pig parts)
 - Conclusion
     - 185 products found containing pig parts
     - Importance of knowing what products are made of
     - Taking better care of raw materials and producers

Figure 5: The detailed process of an iteration in the plan extraction stage

Figure 6: **Topic: This is your brain on communication**. Neuroscientist Uri Hasson researches the basis of human communication, and experiments from his lab reveal that even across different languages, our brains show similar activity, or become äligned,̈when we hear the same idea or story. This amazing neural mechanism allows us to transmit brain patterns, sharing memories and knowledge. Ẅe can communicate because we have a common code that presents meaning,Ḧasson says.

## Initialized Plan

Understanding the Physicality of the Internet
- Architecture and the Internet
  - Physical world vs. digital world
    - Divided attention between screens and surroundings
    - Difficulty in grasping the totality of the Internet
- The Internet's physical reality
  - Squirrel chewing on Internet cable
  - Curiosity about the physical location of the Internet
- The Internet as a real world
  - Visiting data centers and network connection buildings
  - Physical process of connecting networks
- Undersea cables
  - Importance for global connectivity
  - Small in size, expansive in length
  - Light transmission process
- Connecting continents
  - Landing stations and amplifiers
  - Growth of cables in Africa
  - Tata Communications' global belt
- Witnessing a cable connection
  - West Africa Cable System (WACS)
  - Physical process of connecting the cable
- The Internet's history and culture
  - Connection to classic port cities
  - Importance of understanding the physicality of the Internet

## Refined Plan

Understanding the Physicality of the Internet
- Architecture and the Internet
  - Physical world vs. digital world
    - Divided attention between screens and surroundings
    - Difficulty in grasping the totality of the Internet
- The Internet's physical reality
  - Squirrel chewing on Internet cable
  - Curiosity about the physical location of the Internet
  - Physical reality of the Internet
  - Undersea cables connect continents
  - Classic port cities and the internet
- The Internet as a real world
  - Visiting data centers and network connection buildings
  - Physical process of connecting networks
- Undersea cables
  - Importance for global connectivity
  - Small in size, expansive in length
  - Light transmission process
  - Rate of transmission in undersea cables
    - 10-gigabit-per-second wavelength of light
    - Capable of carrying 10,000 video streams
  - Undersea cables use fiber optic technology
- Connecting continents
  - Landing stations and amplifiers
  - Expansion and reliability of Internet connection in Africa
    - Need for stable connection
    - Building an industry around the Internet
    - Ensuring permanent connection
  - Tata Communications' global belt
  - Connection of undersea cables to the continent
- Witnessing a cable connection
  - West Africa Cable System (WACS)
  - Physical process of connecting the cable
  - Fusing undersea cable fibers together
  - Final steps of connecting undersea cables to the shore
- The Internet's history and culture
  - Connection to classic port cities
  - Understanding the infrastructure and physicality of the Internet
  - Losing responsibility for the Internet by thinking of it as a cloud

Figure 7: **Topic: Discover the physical side of the internet**. When a squirrel chewed through a cable and knocked him offline, journalist Andrew Blum started wondering what the Internet was really made of. So he set out to go see it – the underwater cables, secret switches and other physical bits that make up the net.
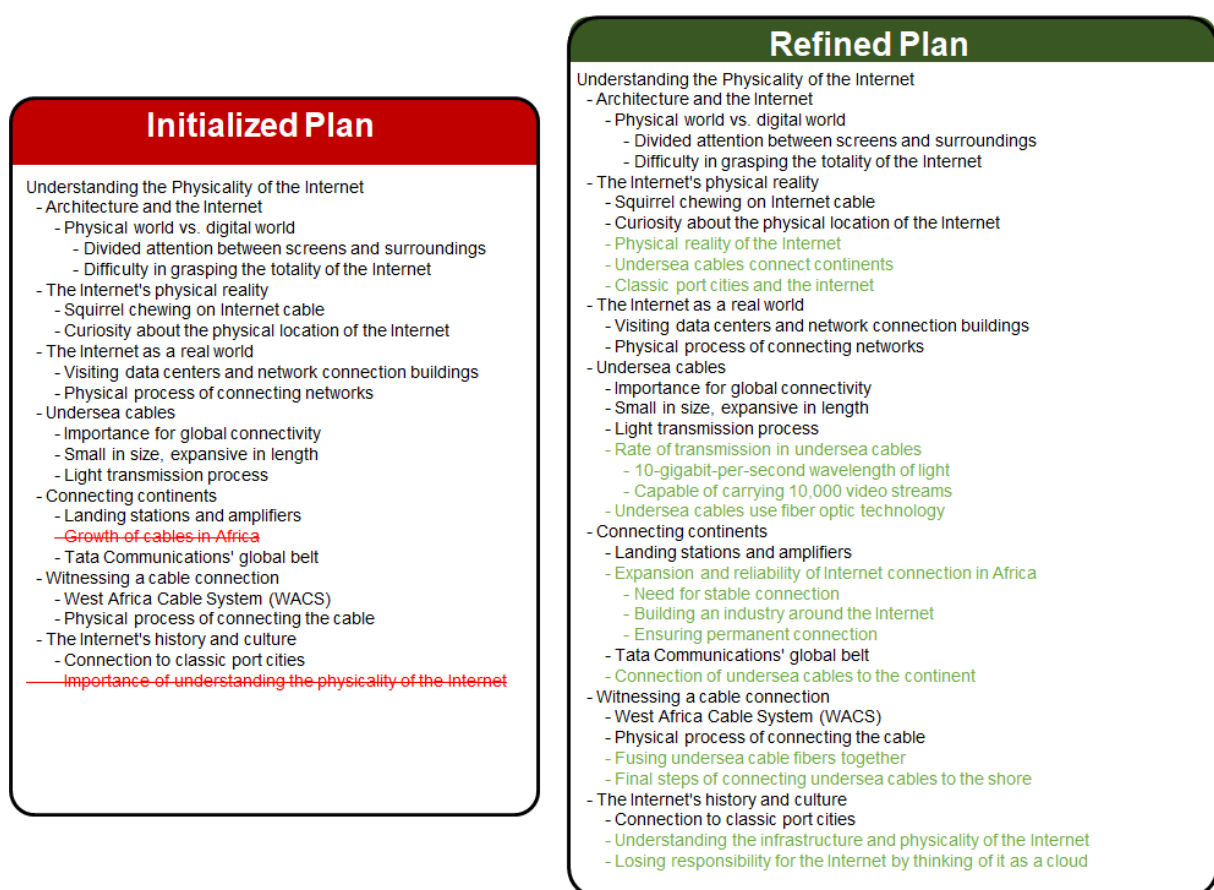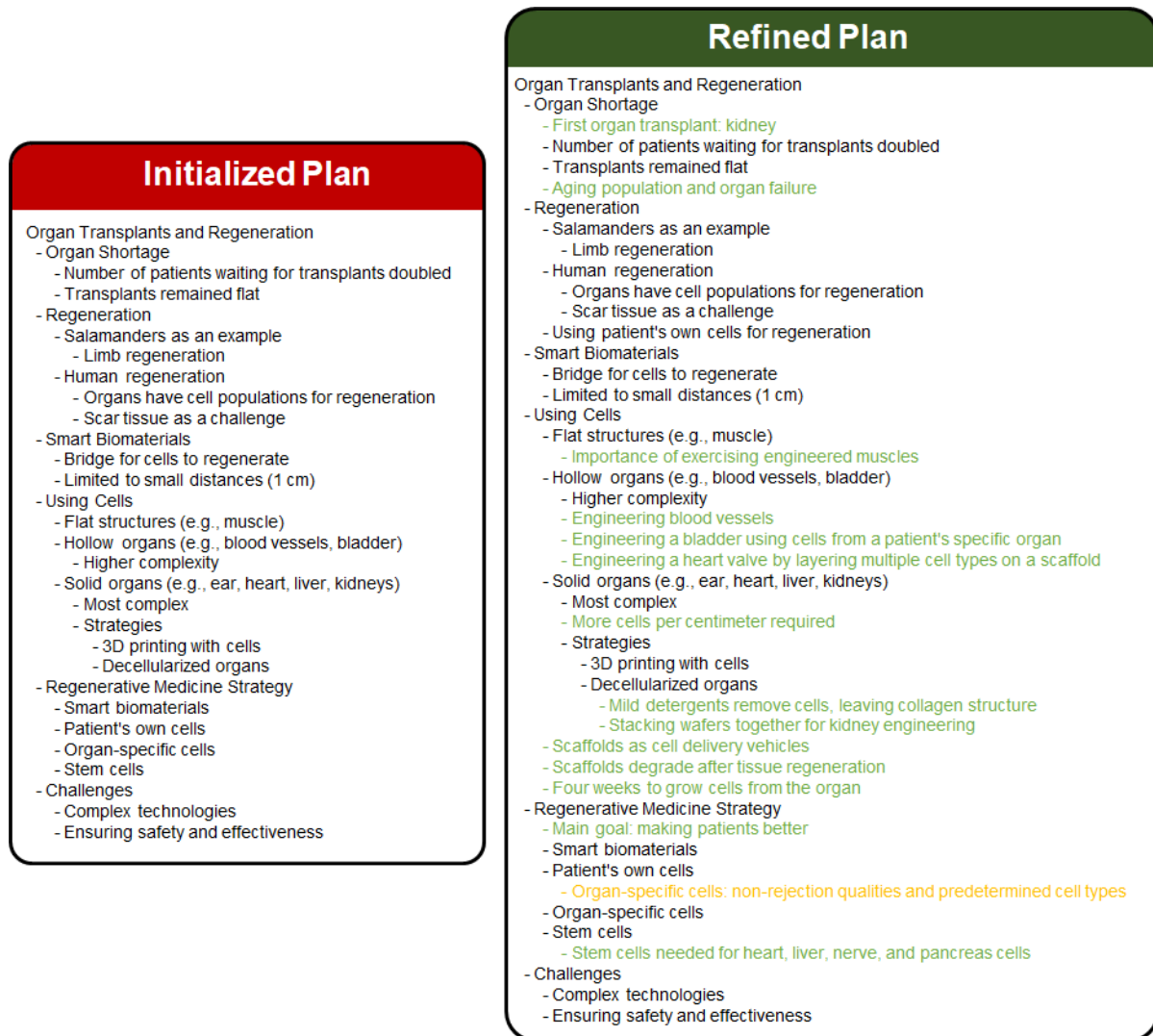
**Initialized Plan**

Organ Transplants and Regeneration
- Organ Shortage
  - Number of patients waiting for transplants doubled
  - Transplants remained flat
- Regeneration
  - Salamanders as an example
    - Limb regeneration
  - Human regeneration
    - Organs have cell populations for regeneration
    - Scar tissue as a challenge
- Smart Biomaterials
  - Bridge for cells to regenerate
  - Limited to small distances (1 cm)
- Using Cells
  - Flat structures (e.g., muscle)
  - Hollow organs (e.g., blood vessels, bladder)
    - Higher complexity
  - Solid organs (e.g., ear, heart, liver, kidneys)
    - Most complex
    - Strategies
      - 3D printing with cells
      - Decellularized organs
- Regenerative Medicine Strategy
  - Smart biomaterials
  - Patient's own cells
  - Organ-specific cells
  - Stem cells
- Challenges
  - Complex technologies
  - Ensuring safety and effectiveness

**Refined Plan**

Organ Transplants and Regeneration
- Organ Shortage
  - First organ transplant: kidney
  - Number of patients waiting for transplants doubled
  - Transplants remained flat
  - Aging population and organ failure
- Regeneration
  - Salamanders as an example
    - Limb regeneration
  - Human regeneration
    - Organs have cell populations for regeneration
    - Scar tissue as a challenge
  - Using patient's own cells for regeneration
- Smart Biomaterials
  - Bridge for cells to regenerate
  - Limited to small distances (1 cm)
- Using Cells
  - Flat structures (e.g., muscle)
    - Importance of exercising engineered muscles
  - Hollow organs (e.g., blood vessels, bladder)
    - Higher complexity
    - Engineering blood vessels
    - Engineering a bladder using cells from a patient's specific organ
    - Engineering a heart valve by layering multiple cell types on a scaffold
  - Solid organs (e.g., ear, heart, liver, kidneys)
    - Most complex
    - More cells per centimeter required
    - Strategies
      - 3D printing with cells
      - Decellularized organs
        - Mild detergents remove cells, leaving collagen structure
        - Stacking wafers together for kidney engineering
  - Scaffolds as cell delivery vehicles
  - Scaffolds degrade after tissue regeneration
  - Four weeks to grow cells from the organ
- Regenerative Medicine Strategy
  - Main goal: making patients better
  - Smart biomaterials
  - Patient's own cells
    - Organ-specific cells: non-rejection qualities and predetermined cell types
  - Organ-specific cells
  - Stem cells
    - Stem cells needed for heart, liver, nerve, and pancreas cells
- Challenges
  - Complex technologies
  - Ensuring safety and effectiveness

Figure 8: **Topic: Growing new organs**. Anthony Atala's state-of-the-art lab grows human organs – from muscles to blood vessels to bladders, and more. At TEDMED, he shows footage of his bio-engineers working with some of its sci-fi gizmos, including an oven-like bioreactor (preheat to 98.6 F) and a machine that p̈rintsḧuman tissue.

**Extracted Plan as Demonstration**

- Powerful Idea about Ideas
  - TED as an example of simplicity
  - Human perception and understanding
    - Shakespeare's Puck: we're easily fooled
    - Betty Edwards: our brain perceives images
      - Faulty perception
      - Seeing objects instead of what's there
      - **Example: two tables with same size and shape**
    - Talmud quote: \"We see things not as they are, but as we are.\"
  - Reality as a hallucination or waking dream
  - Progress through changing perspective
    - Not taught in K-12 curriculum
  - Complexity and simplicity
    - Murray Gell-Mann's emergent properties
    - Gapminder: simplicity without losing important data
    - Cell simulation: missing key aspects of molecular biology
  - Teaching mathematics
    - Pythagorean theorem: different proofs for different levels of understanding
      - Adult sophistication
      - Pythagoras' own proof
    - Six-year-olds discovering mathematical progressions
      - First-order and second-order discrete differential equations
  - Using computers for learning
    - Creating and controlling objects on screen
      - Car example
    - Understanding variables through situated learning
      - Steering wheel example
- Children doing science
  - Experimenting with falling objects
    - Aristotle, St. Thomas Aquinas, and Galileo
    - Importance of experimentation
  - Analyzing acceleration through movie frames
    - Galileo's method: running a ball backwards down the strings of his lute
- Quotes
  - Marshall McLuhan: \"Children are the messages that we send to the future.\"
  - Einstein: \"Things should be as simple as possible, but not simpler.\"
- $100 laptop project
  - Need for mentors and new user interfaces
    - Lack of new teachers
  - Investment in education compared to military spending

**Generated Plan**

The Unexpected Key to Boosting Your Productivity
- The importance of productivity in modern life
- Common misconceptions about productivity
  - Multitasking as a productivity myth
  - The belief that more hours equals more productivity
- The role of the brain in productivity
  - Daniel Kahneman's System 1 and System 2 thinking
    - **Fast, intuitive thinking vs. slow, deliberate thinking**
  - The impact of cognitive load on productivity
    - Example: decision fatigue
- The unexpected key: taking breaks
  - Research on the benefits of taking breaks
    - Improved focus and creativity
    - Reduced stress and burnout
  - Pomodoro Technique: breaking work into intervals
    - Example: 25 minutes of focused work followed by a 5-minute break
- The power of mindfulness and meditation
  - Jon Kabat-Zinn's Mindfulness-Based Stress Reduction (MBSR)
  - Benefits of meditation for productivity
    - Increased focus and attention
    - Enhanced emotional regulation
- Incorporating physical activity into your routine
  - Exercise as a productivity booster
    - Improved cognitive function and memory
    - Reduced stress and anxiety
    - Example: walking meetings or standing desks
- Quotes
  - Peter Drucker: "Efficiency is doing things right; effectiveness is doing the right things."
  - Tim Ferriss: "Focus on being productive instead of busy."
- The importance of a balanced approach to productivity
  - Prioritizing self-care and mental well-being
  - Recognizing the need for rest and recovery
  - Investing in personal development and growth

Figure 9: Example of one shot.