

# Beyond Outlining: Heterogeneous Recursive Planning for Adaptive Long-form Writing with Language Models

Ruibin Xiong<sup>2\*</sup>, Yimeng Chen<sup>1\*†</sup>, Dmitrii Khizbullin<sup>1</sup>, Jürgen Schmidhuber<sup>1,3,4</sup>

<sup>1</sup>Center of Excellence for Generative AI, KAUST <sup>2</sup>Independent Researcher

<sup>3</sup>The Swiss AI Lab, IDSIA-USI/SUPSI <sup>4</sup>NNAISENSE

ruibinxiong@outlook.com, yimeng.chen@kaust.edu.sa

{dmitrii.khizbullin, juergen.schmidhuber}@kaust.edu.sa

<https://www.beyondwriting.org>

## Abstract

Long-form writing agents require flexible integration and interaction across information retrieval, reasoning, and composition. Current approaches rely on predetermined workflows and rigid thinking patterns to generate outlines before writing, resulting in constrained adaptability during writing. In this paper we propose a general agent framework that achieves human-like adaptive writing through recursive task decomposition and dynamic integration of three fundamental task types, i.e. retrieval, reasoning, and composition. Our methodology features: 1) a planning mechanism that interleaves recursive task decomposition and execution, eliminating artificial restrictions on writing workflow; and 2) integration of task types that facilitates heterogeneous task decomposition. Evaluations on both fiction writing and technical report generation show that our method consistently outperforms state-of-the-art approaches across all automatic evaluation metrics, which demonstrate the effectiveness and broad applicability of our proposed framework.

## 1 Introduction

Long-form writing plays a crucial role in numerous domains, including narrative generation (Huot et al., 2024), academic research (Lu et al., 2024), and technical reporting (Shao et al., 2024). Generating coherent, high-quality, and well-structured long-form content presents a significant challenge for Large Language Model (LLM) based writing agents. While LLMs have demonstrated remarkable proficiency in short-form text generation (Yang et al., 2022; Fitria, 2023), their ability to sustain consistency, maintain logical coherence, and adapt dynamically across extended passages remains limited (Yang et al., 2023; Bai et al., 2024; Huot et al., 2024). The complexity of long-form

writing arises from the need to manage interdependent ideas, refine arguments progressively, and integrate diverse information sources, all while ensuring stylistic and factual consistency over extended outputs.

Recent advancements in long-form writing have emphasized a pre-writing planning stage to address these challenges (Yang et al., 2023; Huot et al., 2024; Bai et al., 2024; Shao et al., 2024; Jiang et al., 2024). In the pre-writing phase, an agent first generates a comprehensive outline before proceeding with content generation. For example, Bai et al. (2024) proposed a Plan-Write paradigm to extend LLM-generated content length by planning the structure and target word count for each paragraph then write paragraphs sequentially. Agent’s Room (Huot et al., 2024) argue that a planning stage is important for narrative generation following the narrative theory and proposed a multi-agent framework to generate the plan and write collaboratively. STORM (Shao et al., 2024) incorporates a multi-agent collaborative outlining stage for retrieval-augmented writing.

However, methods that incorporate a pre-writing stage share several limitations. They typically impose a rigid task structure that constrains adaptive reasoning during the writing process, making them less effective for complex, evolving narratives or technical documents that require continual reassessment of writing strategies. Consider a mystery novelist who discovers an unexpected plot element mid-chapter: they must simultaneously retrieve relevant forensic knowledge, reason about plot consistency, and seamlessly integrate new exposition into the narrative flow. Existing structured workflows struggle with such dynamic adjustments because they either establish a fixed outline early in the process or follow a predefined task sequence. This inflexibility prevents writers from making the necessary modifications when they need to revise their plan and engage in deeper reasoning throughout

\*Equal contribution.

†Corresponding author.

the writing process.

In this paper, we unify writing and outlining in a general planning framework. We identify three distinct cognitive tasks involved in the writing process: retrieval, reasoning, and composition, each characterized by unique information flow patterns. Drawing inspiration from Hierarchical Task Network planning (Sacerdoti, 1971; Georgievski and Aiello, 2015), we formulate long-form writing as a planning problem. Within this framework, the overall writing goal is achieved through the execution of primitive retrieval, reasoning, and composition tasks. A valid solution to this problem constitutes a plan that ensures both the executability of each subtask and the achievement of the overarching goal.

Based on the above framework, we propose a general agent structure for long-form writing tasks called heterogeneous recursive planning. Leveraging the goal-directed nature of writing tasks, our approach specifies task types during the planning phase and recursively decomposes them into subtasks across the three cognitive categories: retrieval, reasoning, and composition. This decomposition is recursively applied to subtasks until primitive tasks are reached. The recursive decomposition mechanism enables the system to dynamically adjust planning depth according to the complexity of the writing task. This flexibility allows for efficient adaptation to various document structures and requirements. Incorporating task-type heterogeneity into the planning process facilitates the integration of heterogeneous agents for task execution, while also allowing for type-dependent control mechanisms.

To enable an adaptive writing process, we interleave task execution with planning. When a primitive task is reached, the system immediately executes it, updates the state of all dependent tasks, and then proceeds to the next task node. To manage this execution and recursive planning procedure, we introduce a State-based Hierarchical Task Scheduling algorithm, where tasks and their dependencies are represented as a Directed Acyclic Graph (DAG). We manage the states of tasks to ensuring a hierarchical and dependency-based execution logic.

While existing methods specified to a fixed scenario, such as Agent’s Room for narrative generation and STORM for report writing, we argue that our method can be generalized across multiple writing tasks. We implement our method on two dis-

tinct long-form writing tasks: technical report generation and narrative generation. Our framework is evaluated on relevant benchmarks, including the TELL ME A STORY dataset for fiction writing and the Wildseed dataset for structured document generation. Experimental results demonstrate that our approach significantly improves content quality and adaptability compared to state-of-the-art baselines.

In summary, our key contributions are as follows:

- We propose a planning-based view of long-form writing, casting the process as a set of retrieval, reasoning, and composition tasks. This perspective bridges outlining and writing under a single, goal-driven framework.
- We introduce the heterogeneous recursive planning method that recursively decomposes tasks into subtasks with specified types, enabling flexible integration of specialized agents and supporting type-specific control mechanisms.
- We develop a State-based Hierarchical Task Scheduling algorithm that efficiently manages adaptive execution and dynamic planning.
- Experimental validation of our method demonstrates superior performance over state-of-the-art baselines across diverse writing domains.

## 2 Related Works

**Long-form writing with LLM.** Current approaches to long-form generation primarily adopt a multi-stage paradigm, often designed for specific scenarios with limited generalizability. Early work by Yang et al. (2022, 2023) demonstrates the importance of detailed outlines in story generation. Recently, Bai et al. (2024) proposed that model output length is constrained by SFT data distribution and proposed a Plan-Write framework that extended GPT-4o’s generation to 20,000 words, but maintained a static workflow focused solely on length extension. STORM (Shao et al., 2024), which utilize the autonomous discussion of multi-agents (Li et al., 2023) achieved improved factuality through retrieval-augmented outline generation for Wikipedia-like articles, yet its outlines remain fixed once generated. While Co-STORM (Jiang et al., 2024) further incorporated user interaction for outline optimization in report writing, it still

lacks the capability to dynamically adjust the writing process. Agent’s Room (Huot et al., 2024) employed multi-agent collaboration but imposed rigid role divisions between planning and writing agents, specifically targeting narrative fiction. Although these approaches successfully address their targeted scenarios, their predetermined workflows not only limit adaptability to emergent needs during writing (e.g., contextual conflicts), but also restrict their applicability across different writing tasks.

**Task decomposition.** Neural networks for task decomposition can facilitate long-term sequential planning and decision-making by discovering sub-problems and exploiting sub-solutions (Schmidhuber and Wahnsiedler, 1992). Sec. 5.3 of (Schmidhuber, 2015) describes an adaptive “prompt engineer” which learns to query a separate neural network model for abstract reasoning, planning and decision making. Neural network distillation (Schmidhuber, 1992) can be used to collapse this model and the prompt engineer into a single chain of thought system (Schmidhuber, 2018). Recent research demonstrates that incorporating task decomposition during LLM inference improves performance on language tasks. Wei et al. (2022) showed that explicit chain-of-thought task decomposition during inference significantly enhances the capabilities of LLMs. Approaches like least-to-most prompting (Zhou et al., 2022) and ReAct (Yao et al., 2022) offer improved flexibility in decomposition; however, they are constrained by homogeneous decomposition assumptions, as they treat all subtasks with uniform processing patterns. Khot et al. (2022) designed a planner-executor system with distinct few-shot prompts, but still maintained a rigid operational pattern. Recent approaches like ADaPT (Prasad et al., 2023) introduced on-demand recursive decomposition, yet did not address the integration of fundamentally different types of operations such as retrieval and reasoning. In this work, we propose a heterogeneous recursive framework that effectively handles tasks with distinct operational characteristics.

**Agent workflow.** Agent workflow defines and control the execution logic between sub-modules in an agent system. Several frameworks have been proposed to orchestrate multi-agent workflows. StateFlow (Wu et al., 2024) models workflows as finite state machines, where each task-solving step corresponds to a state with associated output functions, though the methodology for defin-

ing states remains unspecified. While IoA (Chen et al., 2024b)’s Internet-inspired architecture enables multi-device collaboration, it does not address the coordination of cognitive tasks. Recent work has explored search-based optimization of agent workflows (Sordoni et al., 2023; Khattab et al., 2024; Zhuge et al., 2024). For example, AFlow (Zhang et al., 2024) optimizes workflow represented as interconnected action nodes using Monte Carlo Tree Search (MCTS). However task specific optimized workflows remain fixed rather than dynamically adapting to different inputs. This limitation becomes particularly apparent in complex scenarios like long-form writing, where agents need to flexibly alternate between different types of operations based on dynamic context.

### 3 Formulation

In this section, we formulate the fundamental components of a long-form writing agent system, focusing on three heterogeneous task types essential for writing: retrieval (information gathering), reasoning (content planning), and composition (text generation). We further formalize the writing planning problem with a conceptual framework inspired by the hierarchical task network planning.

#### 3.1 Writing Agent System

We first introduce the notion of the writing agent system.

**Definition 3.1** (Writing Agent System). A *writing agent system* is a tuple

$$\Sigma_{\mathcal{A}} = (\mathcal{A}, \mathcal{M}, D, W, \mathcal{O}),$$

where  $\mathcal{A}$  is the *agent kernel* responsible for processing writing instructions, solving writing tasks, and selecting actions.  $\mathcal{M}$  is the *internal memory* maintaining writing-related information like outlines, drafted content, and retrieved references.  $D$  is the *database* (e.g., search engine, reference documents) and  $W$  is the *writing workspace*.  $\mathcal{O}$  is the *input-output interface* through which the agent receives writing instructions and user feedback, accesses information sources, and delivers written content.

#### 3.2 Task Types

The writing process naturally involves three types of heterogeneous cognitive tasks: retrieval for information gathering, reasoning for content planning, and composition for content generation. This

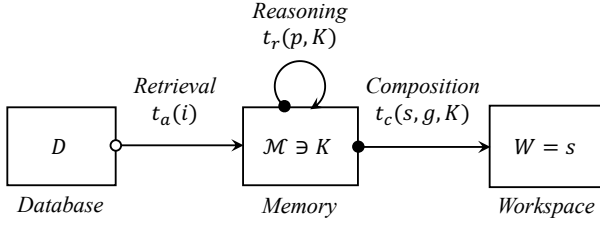


Figure 1: The abstract flow of tasks. The arrow indicates the information flow of a task: the system state at the arrowhead is modified by the labeled task, while the hollow circle end signifies that the associated system state remains unchanged.

categorization aligns with cognitive models of agents (Sumers et al., 2024) and reflects the distinct operational patterns in writing tasks.

**Definition 3.2** (Retrieval Task). Let  $I$  be the set of *information needs* during writing (e.g., factual queries, reference searches). A *retrieval task*  $t_a(i)$  for  $i \in I$  aims to acquire relevant information from the environment and update it into the agent’s memory  $\mathcal{M}$ . This includes searching external knowledge bases, accessing reference documents, or retrieving previously written content.

**Definition 3.3** (Reasoning Task). Let  $P$  represent writing-related problems requiring logical inference (e.g., outline planning, content organization). A *reasoning task*  $t_r(p, K)$  for  $p \in P$  aims to derive new knowledge or make decisions based on available information  $K$  in agent’s internal memory  $\mathcal{M}$ . This includes determining content structure, resolving logical conflicts, or planning writing strategies.

**Definition 3.4** (Composition Task). Let  $G$  represent text generation objectives specifying target states of the written content. A *composition task*  $t_c(s, g, K)$  for  $g \in G$  aims to generate text that meets specified requirements (e.g., style, length, structure) through a sequence of writing actions, given current workspace state  $s$  and knowledge  $K \in \mathcal{M}$ . This includes drafting new content or revising existing text.

We illustrate the abstract flow of the three tasks in Figure 1. Retrieval Task functions as context-independent operations that enhance working memory without modifying the workspace; Reasoning Task performs memory-to-memory transformation contingent upon satisfaction of logical preconditions; and Composition Task executes workspace-altering operations and then updates related infor-

mation to the memory.

### 3.3 Planning for Writing

Planning for writing is based on the assumption that the writing process as complex tasks composed by simpler, executable subtasks. This perspective follows HTN planning, where the objective is not to achieve a set of goals but instead to perform some set of primitive tasks.

In the context of writing, primitive tasks are the basic actions that can be executed directly by the agent. These include retrieval (e.g., searching or accessing documents), reasoning (e.g., analyzing information or ensuring logical consistency), and composition (e.g., generating text, applying revisions, or adding citations). Breaking down complex tasks into these primitives improves accuracy (Chen et al., 2024a) and allows flexible action interleaving.

By assuming a theoretical set  $T_p$  of primitive tasks (without explicitly specifying its composition), we formulate the writing planning problem as follows.

**Definition 3.5** (Writing Planning Problem). A *writing planning problem* is a tuple

$$\langle t_c(g, s_0, K_0), T_p \rangle,$$

where  $t_c(g, s_0, K_0)$  is the top-level composition task, with a writing goal  $g$ , the initial state of the writing workspace  $s_0$ , and the initial content of the agent’s memory  $K_0$ .  $T_p$  is the set of executable primitive retrieval, reasoning and composition tasks.

A solution to this planning problem is a sequence of primitive operations that achieves the writing objective while maintaining coherence and satisfying constraints. Formally, a solution  $\pi = \langle t_1, t_2, \dots, t_k \rangle$  to a writing planning problem must satisfy the following two conditions.

**Condition 1** (Executability). Each operation can be executed when its turn comes, with all prerequisites satisfied.

**Condition 2** (Goal achievement). The final state of  $W$  meets the original writing goal.

## 4 Heterogeneous Recursive Planning

Based on the formulation of the writing task planning problem, we propose a heterogeneous recursive planning method inspired by the HTN planning and the heterogeneity of the three cognitive



tasks. In this section, we introduce the key components of our approach.

#### 4.1 Recursive Planning

The classical HTN planning paradigm solves problems through hierarchical decomposition until reaching primitive executable operations. Following our formulation of the writing planning problem, we adopt a recursive planning strategy, in alignment with classical HTN approaches.

The core of this planning process is task decomposition: each task is broken down into subtasks, and the same decomposition logic is recursively applied to those subtasks. Unlike traditional as-needed decomposition methods that rely on execution failure to stop further planning, our approach introduces a different termination criterion. We only continue planning if certain subtask types necessitate further decomposition, ensuring that the final operations are always executable without redundant decomposition.

#### 4.2 Typed Task Integration

Building upon our formal characterization of cognitive task types in Section 3.2, we extend the recursive planning framework with type-aware decomposition mechanisms.

Our integration addresses the cognitive heterogeneity inherent in writing processes. While complex tasks may involve blended operations (e.g., incidental fact-checking during paragraph composition), their decomposition should respect the dominant cognitive type determined by primary objectives. Our approach is grounded in the following fundamental hypothesis about cognitive task architecture:

**Hypothesis 1** (Type Specification in Decomposition). During hierarchical decomposition of writing tasks, all generated subtasks can be specified as exactly one cognitive type: retrieval, reasoning, or composition.

This hypothesis suggests that the writing planning problem can be decomposed into sub-planning problems of three distinct task types. For example, assume task  $t_c(g, s_0, K_0)$  can be decomposed into a sequential combination of subtasks  $t_a(i)$ ,  $t_r(p, K')$ , and  $t_c(g, s_0, K'')$ , where  $K'$  and  $K''$  denote the modified knowledge in  $\mathcal{M}$  after executing the preceding tasks. The solution of  $\langle t_c(g, s_0, K_0), T_p \rangle$  is then the combination of solutions of planning problems  $\langle t_a(i), T_p \rangle$ ,

$\langle t_r(p, K'), T_p \rangle$ , and  $\langle t_c(g, s_0, K''), T_p \rangle$ . These solutions must satisfy their corresponding executability conditions and goal achievement criteria.

Motivated by the above analysis, we integrate task types into the planning procedure. The benefits of such type specification are twofold. First, it facilitates the function call of heterogeneous agents, for example, a search agent to conduct a retrieval task. Second, it allows different decomposition logic for these inherently different cognitive tasks. For instance, composition tasks require either retrieval or reasoning tasks to modify internal memory, or decompose into multiple subtasks within the writing workspace. They must have a composition-type subtask to reach the goal. Reasoning tasks may decompose into retrieval or reasoning subtasks, either modifying original task conditions or forming logically related subtasks.

In summary, our method features the following design elements:

- **Dynamic type annotation:** Each subtask generated in a planning step is assigned a specific type.
- **Type-aware decomposition:** This provides targeted guidance for potential subtask breakdowns based on the type of the current task.

To ensure goal achievement of the plans, we employ two heuristic decomposition constraints. First, we append a primitive reasoning task after retrieval task sequences to summarize results and update internal knowledge. Second, we require composition tasks to end with a composition subtask to guarantee completion of writing actions in the workspace. These constraints ensure goal attainment through compositional closure in workspace operations and grounding of retrieved information via mandatory reasoning integration.

#### 4.3 Task Scheduling for Adaptive Processing

In our approach, task execution interleaves with planning. Rather than executing subtasks sequentially after completing the entire plan, we execute tasks at any planning node upon termination. This allows the planning and task execution to be adaptive, remaining responsive to previous actions and their outcomes. We present a hierarchical task scheduling algorithm specifically designed for this purpose.

The algorithm addresses dynamic task decomposition depth through sophisticated state man-

---

**Algorithm 1** State-based Hierarchical Task Scheduling

---

**Require:** Directed Acyclic Graph  $G = (V, E)$   
with  $V_{\text{init}} = \{v_{\text{root}}\}$ ,  $S(v_{\text{root}}) \leftarrow \text{ACTIVE}$   
**Ensure:**  $S(v) = \text{SILENCE}$ ,  $\forall v \in V$

- 1: **while**  $\exists v \in V \mid S(v) \neq \text{SILENCE}$  **do**
- 2:   **Find**  $v^* \leftarrow \arg \min_{v \in V} \{\text{BFS-depth}(v) \mid S(v) = \text{ACTIVE}\}$
- 3:   **if**  $v^*$  is primitive task **then**
- 4:     EXECUTETASK( $v^*$ )
- 5:     result  $\leftarrow$  GetExecutionResult()  
      UPDATESATE( $v^*$ ) to SILENCE
- 6:   **else**
- 7:      $\{v_1, \dots, v_k\} \leftarrow \text{TYPEDPLAN}(v^*)$
- 8:     UPDATEGRAPH( $G, \{v_1, \dots, v_k\}, v^*$ )
- 9:     UPDATESATE( $v^*$ ) to SUSPEND
- 10:   **end if**
- 11:   UPDATESATE( $V$ ) to  $\{\text{SILENCE}, \text{SUSPEND or ACTIVE}\}$
- 12: **end while**

---

agement. Tasks and their dependencies are modeled as a Directed Acyclic Graph (DAG). We manage task states (active, suspend, silent) to enforce dependency-based execution order. A task’s state becomes suspended when decomposed into subtasks or when it has unexecuted prerequisites. Tasks become active only when all prerequisites are completed.

Starting from the root task, the algorithm iteratively selects the active task nearest to the root using BFS-based topological sorting. For each selected task, it either executes the task directly if primitive or decomposes it into subtasks that are integrated into the graph. This process continues until all tasks reach the silent state, ensuring systematic traversal and completion of the entire task hierarchy. We summarize the general logic in Algorithm 1.

## 5 Experiments

We evaluate our approach through experiments on challenging language agent tasks addressing three key aspects: (1) the comparative performance of our method against state-of-the-art baselines, (2) the impact of the recursive planning and task-type module, and (3) the generalization capability across diverse task domains. We implemented our proposed method on two long-form writing tasks: narrative generation and report generation.

### 5.1 Narrative Generation

Narrative generation involves complex reasoning and composition tasks. We use the TELL ME A STORY fiction writing dataset proposed in the paper of Agent’s Room (Huot et al., 2024).

**Datasets.** TELL ME A STORY offers a collection of complex, well-structured narratives paired with detailed narrative generation prompts, making it particularly well-suited for evaluating sophisticated story generation capabilities. The dataset consists of 230 samples, with each prompt averaging 113 tokens and corresponding narrative responses averaging 1,498 tokens. The average number of input tokens is 113.

**Baselines.** We evaluate our approach using two state-of-the-art large language models as our base models: Claude-3.5-Sonnet (Anthropic, 2024) and GPT-4o (OpenAI, 2024). For comparison, we implement two primary baselines: (1) End-to-End (E2E): where we directly provide the story prompt to the base LLM without any additional guidance or planning steps; and (2) Agents’ Room (Huot et al., 2024): a collaborative writing framework with multiple agents that decomposes the story generation process into planning and writing phases. In the planning phase, specialized agents outline key story elements including plot structure, character development, and setting details. Writing agents then generate the full narrative following this structured plan.

**Evaluation metrics.** We adopt the LLM-based evaluator for story assessment proposed by Huot et al. (2024), which demonstrates strong correlation with human judgments (Spearman’s rank correlation  $\rho = 0.62, p < 0.01$ ). The evaluation examines four dimensions: plot structure, creativity, character/setting development, and language use. For each story pair, the evaluator determines which is superior or equivalent across these dimensions and overall, producing win-tie-loss judgments. To convert these pairwise comparisons into quantitative scores, we employ the Davidson model (Davidson, 1970), which effectively handles cases with ties. We implement the evaluator using Gemini 2.0 Flash<sup>1</sup> as the base LLM. To mitigate position bias, we conduct  $k$  evaluations in each ordering ( $2k$  total trials) and determine the final outcome through majority voting. We set  $k = 7$  for all results.

---

<sup>1</sup><https://deepmind.google/technologies/gemini/flash/>

Backbones	Methods	Dimensions				
		Plot	Creativity	Development	Language Use	Overall
GPT-4o	E2E	0.337	0.218	0.288	0.202	0.270
	Agent’s Room	1.035	0.712	0.948	0.680	0.869
	Ours	<b>1.470</b>	<b>2.005</b>	<b>1.967</b>	<b>2.233</b>	<b>2.143</b>
	w/o Recursive	1.307	1.327	1.041	1.192	1.100
	w/o Type	0.852	0.733	0.756	0.693	0.717
Claude-3.5-Sonnet	E2E	0.036	0.016	0.032	0.017	0.025
	Agent’s Room	1.029	0.480	0.778	0.484	0.694
	Ours	<b>2.016</b>	<b>2.634</b>	<b>2.959</b>	<b>2.264</b>	<b>2.852</b>
	w/o Recursive	1.145	1.396	0.707	1.517	0.918
	w/o Type	0.774	0.475	0.525	0.518	0.512

Table 1: Quantitative ranking scores of methods across Plot, Creativity, Development, Language Use, and Overall on TELL ME A STORY. Stories generated by methods using the same backbone were compared, with each prompted to generate 2,000 words in length. Ablations of our method are highlighted with a grey background. The highest value in each column is highlighted in bold.

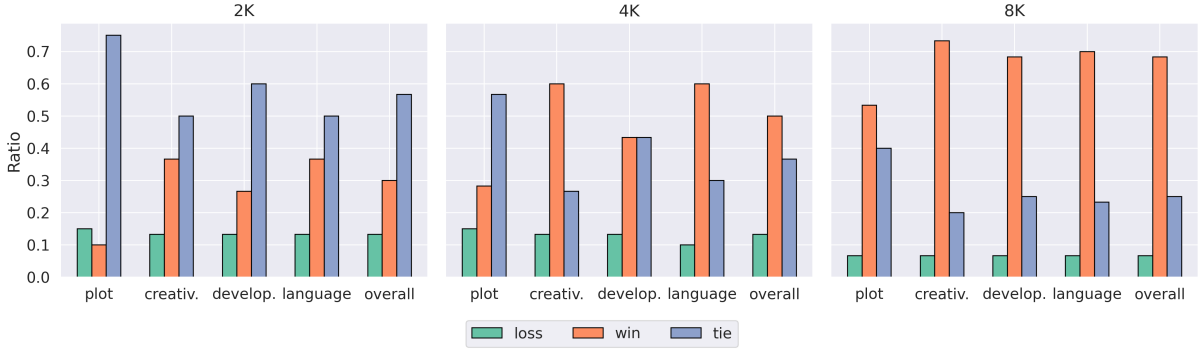


Figure 2: The evaluation results of our method v.s. Agent’s Room at different generation lengths.

**Configurations.** For Agent’s Room baseline, we implement the plan+write version according to the paper, which includes 4 planning agents (conflict, character, setting, plot) and 5 writing agents (exposition, rising action, climax, falling action, resolution). We use a length estimator along with the writing agents to enable the length control. For our method, two task types are included: reasoning (Design) and composition (Writing). We implement a Design agent and a Writing agent as the primitive task executors. For all methods, we employed GPT-4o-20240806 and Claude-3.5-Sonnet-20241022 as the base LLMs with their default parameters.

### 5.1.1 Results

As shown in Table 1, Agent’s Room significantly outperforms the E2E baseline, aligning with results reported in their original paper. Our proposed method demonstrates superior performance across

all five key evaluation metrics compared to baseline approaches. This consistent improvement holds across two different backbone LLMs, validating the robustness of our approach across base models.

**Ablation study.** To analyze the contributions of individual components, we conducted an ablation study with two key variations: 1) Non-recursive generation (“w/o Recursive”): This variant removes the recursive decomposition process, instead generating the entire plan in a single step similar to baseline methods. 2) Task-type removal (“w/o Type”): This variant omits explicit task-type information from the planning stage. While still employing recursive breakdown, the model no longer utilizes type-specific decomposition logic.

**Extended lengths.** We also evaluated how different methods scale with increasing generation length. From our dataset, we selected 60 samples that an LLM identified as suitable for generating

Backbones	Methods	Report Quality			
		Relevance	Breadth	Depth	Novelty
GPT-4o	STORM	4.76	4.58	4.30	4.32
	Co-STORM	4.36	4.22	4.02	4.17
	Ours	<b>4.93</b>	<b>4.86</b>	<b>4.79</b>	<b>4.51</b>
	w/o HRP	4.83	4.18	3.74	4.17
Claude-3.5-Sonnet	STORM	4.66	4.63	4.40	4.41
	Co-STORM	3.87	3.56	3.46	3.82
	Ours	<b>4.96</b>	<b>4.92</b>	<b>4.93</b>	<b>4.82</b>
	w/o HRP	4.84	4.51	4.24	4.46
DeepSeek-R1	Ours	<b>4.97</b>	<b>4.94</b>	<b>4.95</b>	<b>4.88</b>
	w/o HRP	4.94	4.81	4.83	4.80
Commercial	PPL-Deep Research	4.93	4.73	4.75	4.45

Table 2: Comparison of method performance on WildSeek across different backbones under the rubric grading by o1-preview. Our method and its ablations are highlighted with a grey background.

texts over 8,000 words. We then conducted experiments by prompting models to generate articles of three different lengths: 2K, 4K, and 8K words, operating under the assumption that task complexity increases with required text length. Figure 2 presents pairwise comparisons of the overall metric between our method and Agents Room with GPT-4o as the base LLM. We excluded the E2E baseline from this comparison as it is unable to generate texts of 4K or 8K words. For 2,000-word stories, our method and Agents Room performed comparably on more than 50% of samples. However, our method demonstrates increasingly significant advantages over the baseline as task length increases, highlighting its effectiveness in handling more complex long-form content generation.

## 5.2 Report Generation

Compared with story generation, report generation task further need the integration of complex retrieval tasks with reasoning and composition.

**Datasets.** We use the WildSeek dataset proposed by (Jiang et al., 2024). WildSeek offers a collection of real-world information-seeking tasks paired with user goals, making it particularly well-suited for evaluating complex information retrieval capabilities. The dataset consists of 100 samples across 24 domains, collected from 8,777 users of the STORM web application. The data underwent rigorous quality control through rule-based filtering and GPT-4-based classification, followed by manual review. Each data point comprises a Topic-Intent pair, cov-

ering diverse domains including Applied Science, Natural Science, Health & Fitness, Culture & Society, and others.

**Baselines.** We compare our method with STORM (Shao et al., 2024) and Co-STORM (Jiang et al., 2024). STORM is a writing system that uses perspective-guided question asking from retrieval and constructs Wikipedia-like articles through generating outlines and section-by-section writing. Co-STORM extends STORM by introducing a user-participated roundtable discussion to enhance the diversity of retrieved information and improve coverage of unknown unknowns. Both baseline methods rely on retrieval-augmented generation and use similar outline-driven approaches for long-form text generation.

**Evaluation metrics.** We utilize the evaluation framework established by Co-STORM, which examines the final report across four dimensions: Relevance, Broad Coverage (Breadth), Depth, and Novelty. A LLM-based evaluator assesses each dimension on a 5-point scale, with the original Topic and Intent provided. We employ the latest OpenAI o1-preview<sup>2</sup> as our primary evaluator model.

**Configurations.** We use Bing Search API for retrieval. We use the latest official implementation of STORM<sup>3</sup> with their default configurations. For

<sup>2</sup><https://openai.com/index/introducing-openai-o1-preview/>

<sup>3</sup><https://github.com/stanford-oval/storm>



Co-STORM, we follow the official implementation with its user-simulator. Details are provided in Appendix A.3. For both GPT-4o and Claude-3.5-Sonnet we used the same settings as in the narrative generation experiments. We design a search agent, an analyzing agent, and a writing agent as the primitive task executors for retrieval, reasoning and composition respectively. For the search agent, we implement a multi-agent framework comprising a retrieval agent, a reranking agent, and a summarization agent. See Appendix A.4 for more details.

### 5.2.1 Results

Our primary experiment on the WildSeek dataset is presented in Table 2. The results demonstrate that our method consistently outperforms the current state-of-the-art approaches across four distinct automatic evaluation metrics. This further validates the effectiveness and generalizability of our approach. We observe a significant improvement in writing depth with our method, which can be intuitively attributed to the benefits of recursive planning. Additionally, our approach consistently outperforms existing methods in terms of relevance, engagement, and breadth of the generated content.

**Ablation study.** To further validate the effectiveness of our approach, we implemented an ablation version, where we retained the same search agent setup but removed the Hierarchical Recursive Planning method (denoted as “w/o HRP” in Table 2). This means the agent . This allows us to isolate and quantify the performance gains introduced by our proposed approach. We observe a significant improvement in depth metrics compared to the ablation version, demonstrating the impact of our approach. Additionally, removing Hierarchical Recursive Planning results in a notable decline in novelty and breadth, further highlighting its contribution to the generation quality.

**Reasoning model compatibility.** We further experimented using the reasoning model DeepSeek-R1 (DeepSeek, 2024) as the base LLM. Results demonstrate that our approach maintains significant performance advantages. Particularly notable improvements were observed in reasoning depth and breadth metrics. This demonstrates our method’s consistent ability to enhance reasoning capabilities. Our analysis included Perplexity’s Deep Research<sup>4</sup>,

<sup>4</sup><https://www.perplexity.ai/hub/blog/introducing-perplexity-deep-research>

a commercial reasoning model based agent, tested on the same dataset. The results demonstrate that our methodology, when implemented with either Claude-3.5-Sonnet or DeepSeek-R1 as the base model, delivers significantly superior performance across all measured metrics compared to this commercial alternative.

## 6 Conclusions

In this work, we propose a general framework for long-form writing agents, centered around heterogeneous recursive planning. Our approach is based on an analysis of three distinct types of tasks involved in the writing process and formulation of the writing planning problem. We highlight the heterogeneity of writing planning, not only in the final generated plan but also in the sub-planning problems that emerge during hierarchical decomposition. To address this, we incorporate type specification into the recursive planning process. Additionally, we employ a state-based task scheduling algorithm for adaptive task execution. Experiments on narrative generation and report writing demonstrate that our method significantly improves long-form writing quality compared to state-of-the-art baselines. Our ablation studies confirm the importance and effectiveness of both the recursive strategy and task-type awareness, validating their contributions to enhanced content generation.

## 7 Limitations

**Computational efficiency.** The recursive decomposition process introduces additional computational overhead compared to end-to-end approaches. Future work could explore optimization techniques. Another potential avenue for improving efficiency is the use of heterogeneous agents, where models are assigned to different tasks based on their complexity. Instead of applying a single large model to all recursive decomposition and execution steps, specialized models could be leveraged for simpler subtasks, reserving larger models for more complex reasoning.

**Human-in-the-loop integration.** While our approach automates task decomposition and execution, integrating human feedback during the planning and writing stages could further improve adaptability and quality. Exploring interactive refinement mechanisms or reinforcement learning from human evaluations could be promising directions.

## References

- Anthropic. 2024. [Claude 3.5 Sonnet](#). Technical report, Anthropic.
- Yushi Bai, Jiajie Zhang, Xin Lv, Linzhi Zheng, Siqi Zhu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. Longwriter: Unleashing 10,000+ word generation from long context llms. *arXiv preprint arXiv:2408.07055*.
- Qiguang Chen, Libo Qin, Jiaqi Wang, Jinxuan Zhou, and Wanxiang Che. 2024a. Unlocking the boundaries of thought: A reasoning granularity framework to quantify and optimize chain-of-thought. *arXiv e-prints*, pages arXiv–2410.
- Weize Chen, Ziming You, Ran Li, Yitong Guan, Chen Qian, Chenyang Zhao, Cheng Yang, Ruobing Xie, Zhiyuan Liu, and Maosong Sun. 2024b. Internet of agents: Weaving a web of heterogeneous agents for collaborative intelligence. *arXiv preprint arXiv:2407.07061*.
- Roger R Davidson. 1970. On extending the bradley-terry model to accommodate ties in paired comparison experiments. *Journal of the American Statistical Association*, 65(329):317–328.
- DeepSeek. 2024. [DeepSeek-R1: First-rank implementation details of DeepSeek](#). Available at: [github.com/deepseek-ai/DeepSeek-R1](https://github.com/deepseek-ai/DeepSeek-R1).
- Tira Nur Fitria. 2023. Artificial intelligence (ai) technology in openai chatgpt application: A review of chatgpt in writing english essay. In *ELT Forum: Journal of English Language Teaching*, volume 12, pages 44–58.
- Ilche Georgievski and Marco Aiello. 2015. [Htn planning: Overview, comparison, and beyond](#). *Artificial Intelligence*, 222:124–156.
- Fantine Huot, Reinald Kim Amplayo, Jennimaria Palomaki, Alice Shoshana Jakobovits, Elizabeth Clark, and Mirella Lapata. 2024. Agents’ room: Narrative generation through multi-step collaboration. *arXiv preprint arXiv:2410.02603*.
- Yucheng Jiang, Yijia Shao, Dekun Ma, Sina J Semnani, and Monica S Lam. 2024. Into the unknown unknowns: Engaged human learning through participation in language model agent conversations. *arXiv preprint arXiv:2408.15232*.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan A, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2024. [DSPy: Compiling declarative language model calls into state-of-the-art pipelines](#). In *The Twelfth International Conference on Learning Representations*.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406*.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models. *arXiv preprint arXiv:2405.01535*.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. [Camel: Communicative agents for "mind" exploration of large language model society](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 51991–52008. Curran Associates, Inc.
- Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. 2024. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*.
- OpenAI. 2024. [Learning to reason with LLMs](#). Technical report, OpenAI.
- Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal, and Tushar Khot. 2023. Adapt: As-needed decomposition and planning with language models. *arXiv preprint arXiv:2311.05772*.
- Earl D Sacerdoti. 1971. A structure for plans and behavior. *Tech. Note 109*.
- J. Schmidhuber. 1992. Learning complex, extended sequences using the principle of history compression. *Neural Computation*, 4(2):234–242.
- J. Schmidhuber and R. Wahnsiedler. 1992. Planning simple trajectories using neural subgoal generators. In J. A. Meyer, H. L. Roitblat, and S. W. Wilson, editors, *Proc. of the 2nd International Conference on Simulation of Adaptive Behavior*, pages 196–202. MIT Press.
- Jürgen Schmidhuber. 2015. On learning to think: Algorithmic information theory for novel combinations of reinforcement learning controllers and recurrent neural world models. *arXiv preprint arXiv:1511.09249*.
- Jürgen Schmidhuber. 2018. One big net for everything. *arXiv preprint arXiv:1802.08864*.
- Yijia Shao, Yucheng Jiang, Theodore A Kanell, Peter Xu, Omar Khattab, and Monica S Lam. 2024. Assisting in writing wikipedia-like articles from scratch with large language models. *arXiv preprint arXiv:2402.14207*.
- Alessandro Sordani, Xingdi Yuan, Marc-Alexandre Côté, Matheus Pereira, Adam Trischler, Ziang Xiao, Arian Hosseini, Friederike Niedtner, and Nicolas Le Roux. 2023. [Joint prompt optimization of stacked LLMs using variational inference](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

- Theodore Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas Griffiths. 2024. Cognitive architectures for language agents. *Transactions on Machine Learning Research*.
- 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.
- Yiran Wu, Tianwei Yue, Shaokun Zhang, Chi Wang, and Qingyun Wu. 2024. Stateflow: Enhancing llm task-solving through state-driven workflows. *arXiv preprint arXiv:2403.11322*.
- 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. *arXiv preprint arXiv:2210.06774*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, et al. 2024. Aflow: Automating agentic workflow generation. *arXiv preprint arXiv:2410.10762*.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*.
- Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen Schmidhuber. 2024. [GPTSwarm: Language agents as optimizable graphs](#). In *Forty-first International Conference on Machine Learning*.

## A Experiments Details

In this section we introduce the implementation details of our experiments. We also provide additional experiments for the evaluation results. Scores in this section is produced with the open-sourced evaluator LLM Prometheus 2 (Kim et al., 2024) which is shown to have high agreement with the proprietary LM judges.

### A.1 Topic template for WildSeek

In the WildSeek dataset, each sample contains two key fields: Topic and Intent. The Co-STORM paper (Jiang et al., 2024) implements different experimental approaches for these fields. Specifically, for Co-STORM, both the Topic and Intent fields are combined and provided to a LLM that simulates user behavior. In contrast, when using STORM, only the Topic field is supplied to the model. In our implementation we combine the Topic and Intent fields into a refined topic before feeding them to the agents. We remove the trailing period or question mark from the Topic field and make the first letter of the Intent field lowercase. The final refined topic is created with the following template: `f"{topic}, {intent}"`.

### A.2 STORM

Backbone	Method	Tag or branch	Report Quality			
			Relevance	Breadth	Depth	Novelty
GPT-4o	STORM	v1.1.0	4.500	4.530	4.693	4.214
		NAACL-2024-code-backup	4.580 <sup>+0.080</sup>	4.320 <sup>-0.210</sup>	4.617 <sup>-0.076</sup>	3.913 <sup>-0.301</sup>

Table 3: Reproduction experiments for STORM with GPT-4o backbone and Serp/Bing retriever.

Running STORM baseline from the official release branch NAACL-2024-code-backup would involve querying multiple outdated LLMs, specifically, gpt-3.5-turbo, gpt-4, and gpt-4-32k. We use gpt-4o for the fair comparison of the orchestration-level algorithms. We compare the NAACL-2024-code-backup branch to v1.1.0 tag of the official STORM repository in Table 3. We observe that the most recent code v1.1.0 is slightly stronger than the official branch NAACL-2024-code-backup on average across rubrics. We choose tag v1.1.0 as a stronger baseline. We follow the default hyperparameter setting in the official implementation.

### A.3 Co-STORM

Backbone	Method	Input	Variant	Report Quality			
				Relevance	Breadth	Depth	Novelty
GPT-4o	Co-STORM	T	2+I+4-turn	4.429	4.469	4.531	4.255
		T+I	3-turn	4.263 <sup>-0.166</sup>	4.384 <sup>-0.085</sup>	4.535 <sup>+0.004</sup>	3.869 <sup>-0.386</sup>
		T+I	1-turn	4.310 <sup>-0.119</sup>	4.440 <sup>-0.029</sup>	4.380 <sup>-0.151</sup>	4.000 <sup>-0.255</sup>

Table 4: Performance of Co-STORM variants in terms of report quality. The input format ‘T’ refers to that we only include Topic as the original input. ‘T+I’ denotes we combine Topic and Intent as in STORM.

The user collaborative part of Co-STORM is simulated by a LLM. We follow the example implementation provided in the official repository<sup>5</sup>. The user utterance simulation is executed by configuring `costorm_runner.step(simulated_user=True, simulate_user_intent=intent)`. To align with the default setting of STORM, we set `max_search_queries` to 3. We set the number of turns after the

<sup>5</sup>[https://github.com/stanford-oval/storm/blob/main/examples/costorm\\_examples/run\\_costorm\\_gpt.py](https://github.com/stanford-oval/storm/blob/main/examples/costorm_examples/run_costorm_gpt.py)



warm-up phase but before the simulated user utterance to 2, and the number of turns following to 4, thereby simulating 1 and 2 rounds of round-table discussions, respectively.

We conducted a comparative analysis of different Co-STORM variants, as presented in Table 4. The implementation described previously is denoted as 2+I+4-turn. We implemented two variants that use combined Topic and Intent as input, consistent with the approach in both STORM and our method. The two variants differ in the number of turns following the warm-up phase. The first variant employs 3 turns after the warm-up phase without simulated user utterance (designated as "3-turn" in the table). This configuration adheres to the default settings specified in the official example. The other uses just 1 turn (labeled as "1-turn" in the table). The results show that they are relatively worse than the 2+I+4-turn variant, especially in the novelty dimension. We thus present the results of the 2+I+4-turn variant in the main paper.

#### **A.4 Search Agent**

The search agent implementation in STORM and Co-STORM follows a retrieval-augmented generation approach but differs in their information seeking strategies. In STORM, the search agent converts questions into multiple search queries using an LLM, retrieves results through search APIs, and applies rule-based filtering following Wikipedia’s reliable sources guidelines to exclude unreliable sources like social media posts and personal blogs. Co-STORM extends this with a multi-perspective search strategy where agents with different expertise generate questions based on their specialized viewpoints. It also implements a dynamic reranking mechanism that scores retrieved information using a formula which prioritizes information that is relevant to the topic but not directly answering the original question. In our implementation, we employ a multi-agent system as the search agent, consisting of a ReAct-style retrieval agent, a result ranking agent, and a content summarization agent. The retrieval agent issues up to 4 queries and retrieves a maximum of 20 results. These results are then passed to the ranking agent, which scores them and selects the top four based on relevance. The content summarization agent then extracts information from these top-ranked results that is most relevant to the query and search intent, before returning them to the upper-level search execution process.

## B Examples of generated articles

### Fiction

#### Instruction:

Write a story that takes place in a remote cabin in snowy mountains. The story should be written in the first-person perspective. The protagonist has a dog named Sprite who is a whippet. They go to the cabin with Kevin, the protagonist's boyfriend, to celebrate their six-month anniversary. It's their first weekend away together. They go on a hike and get caught in a blizzard. A bear approaches the protagonist and Sprite. There is a plot twist, and it turns out that it is only a dream. The word count should be 5000.

#### Process:

\*1[think].(Design the main characters and their relationships. Include detailed personality traits, backgrounds, and relationship dynamics for the protagonist and Kevin. Also design Sprite's characteristics and role in the story.).

\*2[think].(Design the setting and atmosphere, including detailed descriptions of the cabin, mountain environment, and weather progression. Plan how these elements will contribute to both the real and dream sequences.).1

\*3[think].(Design the plot structure, dream sequence mechanics, and foreshadowing elements that will make the dream revelation both surprising and believable.).1,2

4[write.5000 words].(Write the complete story based on the designed elements, broken into manageable sections.).1,2,3

\*4.1[write.1200 words].(Write the opening section introducing the characters, cabin setting, and initial dynamics. Establish the normal reality while subtly planting seeds for the dream sequence.).

4.2[write.1500 words].(Write the development section showing the first day's activities, building relationship dynamics, and introducing tensions while setting up the dream transition.).

\*4.2.1[think].(Design the evening dinner scene details, including specific romantic moments, initial altitude symptoms, and subtle tensions.).

\*4.2.2[write.500 words].(Write the dinner scene showing the couple's intimacy while introducing subtle tensions and physical discomfort.).4.2.1

\*4.2.3[think].(Design the late evening conversation scene, including discussion topics about Emma's father, Kevin's family history with the cabin, and weather concerns.).4.2.1

\*4.2.4[write.500 words].(Write the evening conversation scene that deepens their connection while introducing underlying concerns.).4.2.3

\*4.2.5[think].(Design the bedtime preparation scene details, focusing on increasing altitude symptoms, weather changes, and Sprite's growing restlessness.).4.2.3

\*4.2.6[write.500 words].(Write the bedtime scene that sets up the dream transition through physical symptoms, weather, and atmosphere.).4.2.5

4.3[write.1300 words].(Write the dream sequence section, including the hike beginning and gradual transition from reality to dream state.).

\*4.3.1[think].(Design the specific mechanics of Emma's transition from sleep to dream state, including physical sensations, environmental changes, and initial dream elements.).

\*4.3.2[think].(Design the initial hike scene details, including weather conditions, character positions, and early warning signs that this might be a dream.).4.3.1

\*4.3.3[write.400 words].(Write the wake-up and transition scene, showing Emma's movement from sleep to what seems like morning, incorporating subtle dream elements.).4.3.1,4.3.2

\*4.3.4[think].(Design the progression of reality distortions during the hike, including changes in environment, time perception, and character behavior.).4.3.2

\*4.3.5[write.450 words].(Write the beginning of the hike scene, establishing the initial seemingly normal situation while introducing subtle dream-like elements.).4.3.4

\*4.3.6[write.450 words].(Write the progression of the hike as reality begins to blur more obviously, leading into the full dream state.).4.3.4

\*4.4[write.500 words].(Write the climactic bear encounter and immediate aftermath.).

\*4.5[write.500 words].(Write the resolution section with the dream revelation and emotional aftermath.).

### **Story:**

The late afternoon sun painted the snow in shades of amber as Kevin's Subaru wound its way up the mountain road. Sprite pressed her elegant whippet nose against the window, her breath creating delicate frost patterns that reminded me of my latest graphic design project. In the driver's seat, Kevin hummed along to some folk song I didn't recognize, his profile relaxed and familiar in the golden light.

"Almost there, Em," he said, reaching over to squeeze my knee. "You doing okay with the altitude?"

I nodded, though the thin air at nearly ten thousand feet made my head feel slightly hollow. The anti-anxiety medication I'd taken earlier wasn't helping the lightheadedness, but I wasn't about to mention that. This weekend meant too much to Kevin—to us both, really. Six months together, and our first real getaway.

The A-frame cabin emerged from between the lodgepole pines like something from a winter postcard. Dark wood and stone stretched upward, with massive windows reflecting the alpenglow. Sprite's tail began its rhythmic whipping against the seat as Kevin pulled into the cleared driveway.

"It's perfect," I breathed, and Kevin's face lit up with that particular smile that still made my stomach flutter. He'd been so excited to share this place, a family property he'd visited since childhood.

The moment I opened the car door, the mountain air hit my lungs like crystallized clarity. Sprite bounded out past me, her sleek form making barely a dent in the fresh powder as she explored the perimeter of the property. The snow crunched satisfyingly under my boots as I made my way to the trunk.

"Let me grab those," Kevin said, already reaching for the heavier bags. "You just get settled inside. The altitude's no joke if you're not used to it."

I wanted to protest—I wasn't some fragile thing—but the way the evening light seemed to wobble at the edges of my vision suggested he might have a point. Instead, I focused on Sprite, who had frozen at the edge of the tree line, her ears pricked forward.

"Sprite? What is it, girl?"

She turned to look at me, then back at the forest. Something about her posture reminded me of that day two years ago when I'd first seen her at the rescue center, all graceful lines and hidden strength. She'd been my anchor during the move to Colorado, during all the nights when memories of Dad's accident had kept me awake.

"Probably just deer," Kevin called out, already heading toward the cabin with our bags. "There's a whole herd that comes through here regularly."

Sprite trotted back to my side, pressing against my leg as we followed Kevin up the wooden steps to the deck. The cabin door opened into a space that felt simultaneously vast and cozy, with a cathedral ceiling that drew the eye upward to exposed beams. A stone fireplace dominated one wall, while the opposite side was nearly all glass, offering a spectacular view of the valley below.

"The sun sets right through those windows," Kevin said, dropping our bags and wrapping his arms around me from behind. "Wait till you see it. Though we might want to get the fire going first—nights get pretty cold up here."

I leaned back against him, watching as Sprite investigated every corner of the space, her nails clicking against the hardwood floors. The walls were decorated with vintage ski equipment and framed photographs of local wildlife—I recognized Kevin's work in several of them, including a stunning shot of a bear silhouetted against a mountain sunrise.

"Your photos look amazing up here," I said, moving closer to examine one of a fox in mid-leap through snow. "The light quality is so different in the mountains."

"Speaking of light—" Kevin checked his watch. "I should get that fire going before we lose the sun. The forecast mentioned a storm system moving in tomorrow."

As if on cue, a gust of wind rattled the windows, making me jump slightly. Kevin pretended not to notice, but I saw the concerned glance he threw my way as he knelt by the fireplace. Sprite appeared at my side again, resting her head against my hip in that intuitive way of hers.

I moved to the kitchen area, unpacking our groceries while Kevin worked on the fire. The space felt both foreign and welcoming, like a place I'd dreamed about but never actually visited. Through the windows, I watched the sun sink lower, casting long shadows across the snow that seemed to reach toward the cabin like grasping fingers.

I shook my head, dispelling the fanciful thought. This weekend was about new beginnings, about Kevin and me taking that next step together. As the first flames caught in the fireplace and Sprite settled onto a plush rug nearby, I felt the familiar knot of anxiety in my chest begin to loosen. We were here. We were safe. Everything else was just shadows and thin mountain air playing tricks with my mind.

The aroma of garlic and herbs filled the cabin as Kevin moved confidently through the kitchen, the wooden spoon dancing between pots with practiced ease. I sat at the counter, trying to focus on chopping vegetables while the room performed a lazy spin around me.

"You've got to try this," Kevin said, bringing a spoon to my lips. His family's Irish stew recipe, he'd explained earlier, handed down through generations of Walsh winter nights. The broth tasted oddly metallic—the altitude playing tricks with my senses—but I smiled at his expectant look.

"Perfect," I lied, the word slightly breathy. Even sitting, the thin air made every movement feel like swimming through silk.

Sprite's nails clicked against the hardwood as she moved from her spot by the fire to press against my legs. Kevin's hand found my shoulder, thumb tracing small circles. "Why don't you let me finish up here?"

"I'm fine," I said, perhaps too quickly. The knife slipped, nearly catching my finger. "Just need to concentrate."

Kevin's lips brushed my temple as he reached past me for the herbs, his chest warm against my back. "You know," he murmured, "Mom always said altitude makes everything taste different. Maybe that's why her recipes never worked in Denver."

The wind picked up outside, rattling the massive windows. Sprite's ears twitched at the sound, her eyes fixed on the darkening glass. The snow was falling faster now, thick flakes swirling in patterns that made my vision swim.

We swayed together by the stove, Kevin humming that same unfamiliar folk song from the car while stirring the stew. My head rested against his shoulder, as much for stability as intimacy. The room felt too warm, then too cold, my body unable to decide which.

"Remember our first date?" Kevin asked, adding another pinch of salt. "When you ordered that ridiculously spicy curry to impress me?"

"And ended up drinking a gallon of water?" The memory brought a genuine laugh, though it left me slightly breathless. "At least I made an impression."

"You always do." He turned, catching my waist as I swayed. "Steady there, Em."

"Just moved too fast," I said, but allowed him to guide me to a chair. Sprite immediately rested her head in my lap, her presence grounding.



The lights flickered once as Kevin served the stew into deep bowls. Through the windows, the snow had transformed the world into a white blur, the trees barely visible. Tomorrow's hike hovered in my thoughts, unspoken between spoonfuls of stew that I couldn't quite taste properly.

"Storm's coming in faster than they predicted," Kevin remarked, his tone carefully casual. "Good thing we're staying in tomorrow."

I stirred my stew, watching the vegetables swirl like the snow outside. "We are still hiking tomorrow, right?" The words came out steady, practiced.

His hesitation lasted only a heartbeat, but it echoed in the space between us like the wind's hollow moan.

After dinner, we migrated to the leather couch facing the fireplace, the empty bowls abandoned on the coffee table. The storm pressed against the windows like a living thing, making the flames dance and flicker. Sprite curled between us, her slender body radiating warmth.

"You know," Kevin said, his fingers absently tracing patterns on my shoulder, "my grandfather used to say this cabin had a way of bringing out truths in people. Something about the isolation, maybe, or how the mountains strip everything down to essentials."

I watched the fire cast his face in amber light and shadow. "Is that why you brought me here?"

"Partly." He shifted, reaching for something beside the couch—a leather-bound album I hadn't noticed before. "I've never brought anyone else up here. It's always been just family."

The album's pages crackled as he opened them, revealing faded photographs of a younger cabin, its wood still raw and new. A man with Kevin's eyes and smile stood proudly in front of it, arm around a woman in a vintage ski jacket.

"Grandpa built it himself in '65," Kevin said. "Said he'd planned to build it elsewhere, but when he found this clearing, the mountains told him this was the spot." His voice softened. "He died up here, you know. Heart attack while photographing a sunrise. Mom says it's exactly how he would have wanted to go."

The wind howled a counterpoint to his words. Sprite's ears twitched, and she pressed closer to my leg.

"My dad used to say mountains show you who you really are," I found myself saying, the words rising unbidden. "That they don't care about your plans or preparations. They just are, and you have to meet them on their terms."

Kevin's hand found mine in the firelight. "Is that what worries you about tomorrow? Meeting the mountain on its terms?"

I stared into the flames, remembering another fire, another night. "The last time Dad went hiking, he had all the right gear, knew all the right moves. The mountain didn't care."

"Em..." Kevin's voice was gentle, but I could hear the familiar tension beneath it—the careful balance between pushing and protecting.

A particularly fierce gust rattled the windows, making us both jump. The weather alert on Kevin's phone chirped, casting a brief blue glow over our faces.

"Storm's intensifying faster than expected," he said, studying the screen. "Maybe we should—"

"I want to do the hike," I interrupted, the words coming out firmer than I felt. "I need to."

Kevin was quiet for a long moment, his thumb tracing my knuckles. Finally, he nodded. "Okay. But we do it smart. We do it safe."

Sprite lifted her head suddenly, staring at the dark windows with an intensity that made my skin prickle. Beyond the glass, the swirling snow seemed to form shapes that dissolved as quickly as they appeared, like memories slipping through my fingers.

The stairs to the loft seemed steeper than they should be, each step requiring more concentration than the last. Kevin's hand at the small of my back steadied me, but the touch felt distant, as if coming through layers of cotton. Above us, the skylight framed a dizzying dance of snowflakes that made the room tilt slightly.

"Easy there," Kevin murmured as I stumbled on the final step. "The altitude really hits you up here."

Sprite darted past us, her usual fluid grace replaced by an anxious energy. She paced the perimeter of the loft, pausing at each window to stare into the whirling darkness. The wind found new voices in the cabin's bones, whistling through invisible gaps with an almost musical persistence.

My fingers fumbled with pajama buttons, the simple task made complex by altitude-numbed hands and the strange way the shadows moved across the walls. Kevin stepped in to help, turning it into an intimate moment, but Sprite's sudden growl at the skylight shattered the warmth of his touch.

"Just the storm, girl," Kevin said, but his voice carried an edge I'd never heard before.

In the bathroom, I gripped the sink's edge as the room performed a lazy spin. The metallic taste in my mouth intensified, and my reflection in the mirror seemed to lag slightly behind my movements. Through the door, I heard Sprite's nails clicking an erratic rhythm on the hardwood, punctuated by soft whines.

"Here," Kevin appeared with water and aspirin, his form wavering slightly in the doorway. "Mom always said altitude sickness gets worse at night."

The storm surged against the cabin, and the lights flickered once, twice. In that stuttering darkness, the shadows cast by the skylight seemed to move independently of their sources, reaching across the ceiling like grasping fingers.

Sprite refused to settle in her usual spot at the foot of the bed, instead pressing herself against my legs, her body trembling with an energy I could feel through the mattress. The digital clock on the nightstand blinked numbers that didn't quite make sense, and I couldn't remember if I'd taken my evening medication or just thought about taking it.

"The storm's really picking up," Kevin said, his voice sounding simultaneously close and very far away. He checked his phone, frowning. "Coverage is getting spotty, but the alert says—"

The wind drowned out his words, its howl transforming into something almost vocal. Sprite's head snapped toward the skylight, tracking something I couldn't quite see in the swirling snow. The room felt too warm, then abruptly too cold, and the ceiling seemed to breathe with the storm's rhythm.

"Sleep," Kevin whispered, pulling me close. "Everything looks better in the morning."

But as I drifted off, Sprite's low growl vibrated through the mattress, and the shadows continued their silent dance across the walls, telling stories I wasn't sure I wanted to understand.

Consciousness returned in layers, each one more uncertain than the last. The weight of Sprite against my legs felt both present and impossibly distant, as if she existed in two places at once. Through the skylight, snow was still falling—or had it stopped and started again? The quality of light seemed wrong somehow, too golden for dawn but too dim for afternoon.

"Kevin?" My voice sounded hollow in my ears, echoing slightly more than the room's acoustics should allow. The space beside me was empty, the sheets cool as if he'd been gone for hours, yet I could smell his coffee as clearly as if he were standing next to the bed.

Sprite's warmth disappeared from my legs, and I heard her nails on the hardwood—click-click-click-click—but the rhythm seemed to continue long after she'd stopped moving. When I pushed myself up, the room tilted at impossible angles before settling into something that almost resembled normal geometry. The digital clock on the nightstand blinked 8:47, then 10:23, then 8:47 again.

"Just the altitude," I murmured, but the words tasted like metal and pine needles.

The morning light through the windows cast shadows that moved independently of the swaying trees outside. I watched, transfixed, as they crawled across the floor like living things, forming patterns that almost resembled footprints in snow. Somewhere below, I heard Kevin's laugh, followed by the clink of coffee cups, but the sound seemed to come from multiple directions at once.

Sprite appeared in the doorway, her elegant form backlit by impossible sunlight. Her eyes caught the light and reflected it back with an intensity that made my head swim. She whined—a sound that started normal but stretched into something musical and strange—then turned and disappeared down the stairs.

My feet found the floor, which felt simultaneously solid and slightly fluid, like walking on packed snow that hadn't quite decided to melt. The air grew thicker with each step toward the stairs, carrying scents that shouldn't go together: coffee, pine, Kevin's aftershave, and underneath it all, the metallic tang of approaching snow.

The morning stretched like pulled taffy as we prepared for the hike, time seeming to catch and release in strange intervals. Kevin laid out our gear with methodical precision—each item appearing on the cabin's wooden floor in perfect alignment, though I couldn't quite remember watching him place them there.

"Trail should be packed down enough for these," he said, holding up my snowshoes, though his voice seemed to come from somewhere slightly left of where he stood. "Three miles up to Thompson's Ridge, then back before the storm hits."

I nodded, swallowing another altitude pill that dissolved with an electric tingle on my tongue. Through the window, the sun hung like a pale coin in a sky that couldn't decide between blue and white, casting shadows that moved a fraction too slowly across the snow.

Sprite paced circles around our gear, her usual pre-walk excitement transformed into something more urgent. Her paws left prints in the cabin's hardwood that seemed to linger a moment too long before fading, like afterimages from staring at the sun. When she paused to stare out the window, her reflection showed three distinct silhouettes before merging back into one.

"Ready?" Kevin's hand appeared on my shoulder, warm through layers that felt simultaneously too thick and too thin. The air inside the cabin had developed a crystalline quality, refracting morning light into prisms that caught in my peripheral vision.

Outside, the snow crunched beneath our boots with a sound that echoed slightly out of sync with our steps. Sprite bounded ahead, her sleek form moving with impossible grace through snowdrifts that seemed to shift and reshape themselves in her wake. The trail marker read "Thompson's Ridge - 3.2 miles," though I could have sworn it had said 2.8 when Kevin first pointed it out.

"Stay close," Kevin called, his figure already seeming somehow less substantial against the white backdrop of snow and sky. "Storm's probably moving faster than the forecast showed."

I checked my watch—9:47 AM—then again—10:22 AM—then once more—9:47 AM. The thin mountain air crystallized in my lungs, each breath carrying the taste of approaching snow and something else, something metallic and familiar that I couldn't quite name. Sprite returned to my side, her warm presence the only constant in a landscape that seemed to be slowly, subtly, rearranging itself around us.

The trail ahead split into three identical paths, then merged back into one as I blinked. Kevin's figure wavered like heat rising from summer pavement, though the air bit cold enough to crystallize thoughts. Sprite darted between trees that seemed to shift positions when I wasn't looking directly at them, her whippet form leaving prints in the snow that filled with impossible colors.

"The ridge should be just ahead," Kevin's voice echoed from multiple directions, though he stood right beside me. Or had he moved ahead? The snow-laden branches above us cast shadows that moved against the wind, reaching down like grasping hands.

My watch face swirled with numbers that refused to settle. 11:47. 10:13. Yesterday. Tomorrow. Time stretched like taffy, then snapped back with a force that left me gasping. The metallic taste in my mouth intensified, familiar as the copper penny scent of Dad's climbing gear.

"Em?" Kevin called, his voice distorting. "You okay back there?"

I tried to respond, but the words froze in the air, hanging like icy crystals before shattering. Sprite pressed against my legs—once, twice, three times simultaneously—her usually sleek form rippling with impossible grace. The trail markers we passed told different stories: Summit 2.4

miles. Base 5.7 miles. Home was never here.

The storm rolled in like a living thing, snow falling upward in geometric patterns that wrote equations in the air. Kevin's red jacket multiplied in the whiteness—ahead of me, beside me, behind me—each version slightly different, slightly wrong. The mountain itself seemed to breathe, expanding and contracting with each step we took.

"We should turn back," I heard myself say, but the words came out in my father's voice. Sprite's ears pricked forward, tracking sounds that existed somewhere between memory and prophecy. The snow beneath our feet had become transparent, showing other trails, other hikers, other times layered like geological strata.

Kevin turned—or had he been facing me all along?—his features blurring at the edges. "The cabin's closer than home," he said, though I couldn't remember which direction either lay. The wind carried the scent of woodsmoke and tomorrow's breakfast, mixed with the sharp tang of fear.

Sprite froze, her elegant form suddenly too still, too perfect, like a photograph of herself. Through the swirling snow, dark shapes moved with deliberate purpose, and I realized the trees had been walking alongside us all along, their branches reaching, reaching—

The world tilted sideways, reality peeling away like old wallpaper to reveal the dream beneath.

Through the crystalline chaos of snow, the bear emerged like a shadow gaining substance. Its form seemed to absorb the whirling fractals of ice, growing larger with each step, its edges bleeding into the white void where reality had been. Sprite's growl vibrated at an impossible frequency, her slender form elongating as she placed herself between me and the approaching mass of midnight fur and memory.

"Dad?" The word escaped in a cloud of frozen breath that shattered into prisms. The bear's eyes held the same amber warmth as the cabin's windows, as Kevin's morning coffee, as the last sunset I'd watched with my father before—

Sprite launched herself at the bear, her whippet form stretching like mercury, multiplying into a dozen silver arrows that pierced the space between heartbeats. The bear rose, and rose, and rose, its shadow falling upward into the inverse snow. Kevin's voice echoed from everywhere and nowhere: "The mountains show you who you really are."

My father's climbing rope unspooled from the bear's mouth, crystallized with decades of ice, spelling out words in its coils: trust, fear, fall, catch. Sprite's barks came in perfect thirds, harmonizing with the wind's hollow chorus. The bear's paw—massive as mortality, gentle as Kevin's morning touch—reached not for me, but through me, pulling out all the tangled knots of anxiety I'd carried up this impossible mountain.

"I can't lose anyone else," I whispered, but the words transformed into snowflakes, each one containing a tiny scene: Kevin leading me up the trail, Dad checking his harness one last time, Sprite pressing against my legs during midnight panic attacks. The bear gathered these moments in its massive paws, holding them up to the fractal sky.

Sprite's light blazed suddenly, her protective love burning away the storm's edges, revealing the seams where dream met memory. The bear began to dissolve, each particle of its being a star falling upward, forming constellations of all my unspoken fears. As it faded, its eyes met mine with infinite understanding—they were Kevin's eyes, they were Dad's eyes, they were my own reflected in the cabin's windows at dusk.

The world tilted once more, reality rushing back like an indrawn breath. My eyes snapped open to find Sprite's warm weight across my legs, solid and real in the cabin's darkness. Kevin's arm was draped protectively over my waist, his breathing deep and steady. Through the skylight, stars wheeled in familiar patterns, and somewhere in the distance, a real bear ambled through real snow, both more and less magnificent than its dream counterpart.

I pressed my hand against my chest, feeling my heartbeat settle into the same rhythm as the mountain's ancient pulse. Tomorrow, we would hike the real trail, face the real weather, and I would carry both the dream's wisdom and its warning. For now, though, I simply breathed, letting



the thin mountain air fill my lungs with possibility rather than fear.

The cabin's warmth enveloped me as consciousness settled like fresh snow. Sprite's steady breathing against my legs anchored me to reality, her slender form no longer mercury-quick and impossible. Through the skylight, dawn painted the ceiling in watercolor washes of pink and gold, all the dream's fractal patterns resolved into simple morning light.

"Hey." Kevin's voice, sleep-rough but present, came from beside me. "You were dreaming pretty intensely there."

I turned to find his eyes—just his eyes, not the bear's, not my father's—watching me with gentle concern. The mountain air still felt thin, but now it carried only the clean scent of pine and possibility, the metallic tang of fear washed away by morning.

"I was—" The words caught, not from altitude this time, but from the lingering echo of dream-truth. "I've been so afraid of losing this. Of losing you, like I lost him."

Sprite shifted, pressing her warm length more firmly against me, and Kevin's hand found mine in the quiet space between heartbeats. Outside, the storm had passed, leaving behind a world crystallized in new snow and clarity.

"The mountains don't take things from us, Em," he said softly, echoing a conversation we'd had a lifetime ago at that first farmer's market. "They just show us what we're carrying."

I thought of the dream-bear, gathering my fears in its massive paws like precious stones. "I think I've been carrying him with me—Dad, I mean—but not in the way I should have been."

Kevin sat up, his shoulder solid against mine, and we watched the sunrise paint the valley in shades of possibility. Sprite's tail thumped a gentle rhythm against the mattress, her dark eyes holding none of the dream's impossible geometries.

"You know what else these mountains show us?" Kevin's voice held the same warmth as the cabin's morning-lit windows. "They show us who's still here, still climbing beside us."

The truth of it settled in my chest like good whiskey, warm and clarifying. I had been so focused on who might leave that I'd nearly missed who had stayed—Kevin with his patient strength, Sprite with her graceful courage, and me, finally ready to trust the trail ahead.

"So," I said, watching a pair of ravens trace figure-eights against the crystal sky, "about that hike?"

Kevin's smile broke like dawn. "Whenever you're ready. The mountain's not going anywhere."

And neither, I realized, was I. The dream's wisdom had crystallized into something I could carry, lighter than fear and stronger than memory. Outside our window, the real world waited, full of real challenges and real joys, all of them possible, all of them mine to choose.

## Report

### Instruction:

Write an article about DeepSeek R1, focusing on the technology it represents and its impact on the future of the large model industry.

### Report:

\* DeepSeek R1: Revolutionizing the Large Model Industry

\*\* I. Introduction: The Emergence of DeepSeek R1

In January 2025, the artificial intelligence landscape witnessed a potential inflection point with the release of DeepSeek R1, a groundbreaking large language model developed by the Chinese AI startup DeepSeek [11][14][16]. While new AI models emerge with increasing frequency in today's rapidly evolving technological environment, DeepSeek R1 stands apart as more than just another incremental advancement—it represents a fundamental rethinking of how large language models

are designed, deployed, and democratized.

At first glance, DeepSeek R1's headline specifications appear impressive but not revolutionary: a 671 billion parameter model with a 128K token context window [1][2][5]. However, beneath these numbers lies an architectural innovation that challenges the industry's "bigger is better" paradigm. Unlike traditional models that activate all parameters for each token processed, DeepSeek R1 selectively activates just 37 billion parameters per token—approximately 5.5% of its total parameter count [1][2][3]. This parameter efficiency enables DeepSeek R1 to achieve performance comparable to models with far greater computational requirements, including OpenAI's o1 series [4][7][8].

Perhaps equally significant is DeepSeek's decision to release R1 under the MIT License, making it fully open-source for commercial use, modifications, and derivative works [1][2][5][9]. This approach stands in stark contrast to the increasingly proprietary nature of leading AI models and signals a commitment to democratizing access to cutting-edge AI capabilities. By combining state-of-the-art performance with open accessibility, DeepSeek R1 challenges the notion that advanced AI must remain the exclusive domain of well-resourced technology giants.

The model's cost profile further reinforces its disruptive potential. With estimated training costs of approximately \$5.58 million—a fraction of the reported \$100+ million for models like GPT-4—and operational costs around \$8 per million tokens, DeepSeek R1 dramatically lowers the financial barriers to deploying advanced AI [3][9][14]. This cost efficiency could accelerate AI adoption across industries and geographies, particularly in emerging markets and among smaller organizations previously priced out of the advanced AI ecosystem.

DeepSeek R1's emergence also carries significant geopolitical implications. As a Chinese AI startup achieving performance comparable to leading Western models, DeepSeek represents a potential rebalancing of global AI capabilities [20]. This development occurs against a backdrop of intensifying technological competition between the United States and China, with both nations heavily investing in AI as a strategic priority [20].

Beyond its technical specifications and market positioning, DeepSeek R1 has demonstrated exceptional capabilities in complex reasoning tasks. Its performance on mathematical benchmarks like AIME 2024 (79.8%) and MATH-500 (97.3%) exceeds even OpenAI's o1-1217 model in certain areas [4][7][8]. This mathematical prowess, combined with strong performance across coding, multilingual understanding, and general knowledge tasks, positions DeepSeek R1 as a versatile foundation for diverse applications.

In the following sections, we will explore DeepSeek R1's technical foundations, performance profile, market positioning, and potential impact on the large model industry. We'll examine how its innovative architecture, training methodology, and deployment approach could reshape industry standards and practices. Finally, we'll consider the future development directions and challenges that may emerge as DeepSeek R1 and similar models continue to evolve, potentially marking a new phase in AI development characterized by efficiency, accessibility, and specialized capabilities rather than simply scaling up existing approaches.

## **\*\* II. Technical Foundation: Architecture and Innovations**

At the heart of DeepSeek R1's groundbreaking capabilities lies a sophisticated technical architecture that challenges conventional approaches to large language model design. While the AI industry has historically pursued performance gains through ever-larger models, DeepSeek R1 represents a fundamental shift toward parameter efficiency, selective computation, and specialized training methodologies. This section examines the four key technical innovations that collectively position DeepSeek R1 as a potential inflection point in large model development.

### **\*\*\* Parameter Efficiency Revolution**

DeepSeek R1's most striking innovation is its approach to parameter efficiency. The model contains a massive 671 billion total parameters, yet activates only 37 billion parameters (approximately 5.5% of the total) for processing each token [1][2][3]. This selective activation strategy dramati-

cally reduces computational requirements during both training and inference while maintaining performance comparable to models that activate all parameters for every token [3][7].

To put this efficiency in perspective, DeepSeek R1 achieves similar or superior performance to OpenAI's o1 series while activating approximately 50 times fewer parameters per token than GPT-4's reported 1.8 trillion active parameters [3]. This parameter efficiency translates directly to lower computational costs, reduced energy consumption, and greater accessibility—addressing three of the most significant challenges facing large-scale AI deployment [3][9].

The parameter efficiency revolution embodied by DeepSeek R1 represents more than just an incremental improvement; it suggests a fundamental rethinking of how large language models should be designed. Rather than simply scaling up existing architectures, DeepSeek has prioritized selective computation through a sophisticated routing mechanism that determines which parameters should be activated for each specific input [3][7][9].

#### \*\*\* Mixture of Experts Framework

The technical foundation enabling DeepSeek R1's parameter efficiency is its implementation of the Mixture of Experts (MoE) framework [3][7][9][10][15]. Unlike traditional transformer models that process all inputs through the same parameters, MoE architectures contain multiple specialized neural networks called "experts," with a routing mechanism that selectively activates only the most relevant experts for each token [3][7][15].

DeepSeek R1's implementation includes several advanced MoE techniques that enhance its efficiency and performance:

1. **Auxiliary-Loss-Free Load Balancing**: The model implements a dynamic bias term for routing updates that doesn't influence model weight updates, ensuring balanced utilization of experts without compromising performance [15].
2. **Shared Experts**: DeepSeek R1 incorporates shared experts that learn from all tokens, enhancing the model's ability to generalize across diverse inputs while maintaining specialization where beneficial [15].

This sophisticated MoE implementation enables DeepSeek R1 to achieve the seemingly contradictory goals of massive scale (671B total parameters) and computational efficiency (37B activated parameters per token), representing a significant advancement over traditional approaches to large model design [3][7][9].

#### \*\*\* Advanced Attention Mechanisms

DeepSeek R1 incorporates Multi-Head Latent Attention (MLA), a sophisticated attention mechanism that further enhances computational efficiency while maintaining the model's ability to process its extensive 128K token context window [10][15]. MLA compresses Query, Key, and Value vectors to optimize the computational resources required for each token processed, addressing one of the most resource-intensive components of transformer architectures [10][15].

This attention optimization is enhanced by several complementary techniques:

1. **Dynamic Low-Rank Projection**: The model adaptively adjusts compression based on sequence length, optimizing for both short queries and long-context reasoning tasks [10][15].
2. **Adaptive Query Compression**: DeepSeek R1 varies query scaling at different model depths, enhancing the model's ability to capture both fine-grained details and broader contextual patterns [10][15].

3. **Joint KV Storage**: The model implements efficient memory management techniques that reduce memory traffic during inference, further enhancing computational efficiency [10][15].

These attention mechanisms collectively enable DeepSeek R1 to process and reason over extremely long contexts while maintaining computational efficiency—a critical capability for complex reasoning tasks that require integrating information across thousands of tokens [1][2][5][11].

#### \*\*\* Training Methodology Breakthroughs

DeepSeek R1's capabilities are further enhanced by its innovative multi-phase training methodology, which combines supervised fine-tuning with sophisticated reinforcement learning techniques

[2][4][5][7]. This hybrid approach consists of four distinct phases:

1. **Cold Start Phase** (2 weeks): Initial fine-tuning on a carefully curated dataset of chain-of-thought examples to establish foundational reasoning capabilities [3][7][8][10][14].
2. **Reinforcement Learning Phase** (8 weeks): Large-scale reinforcement learning across diverse tasks to enhance the model’s reasoning abilities and alignment with human preferences [3][7][8].
3. **Rejection Sampling and Fine-Tuning Phase**: Selection of optimal responses from multiple generated outputs for further training refinement, enhancing response quality and consistency [7][10].
4. **Diverse Reinforcement Learning Phase**: Training across varied tasks with rule-based rewards to strengthen specific reasoning capabilities and address weaknesses identified in earlier phases [7][8].

A particularly notable innovation in DeepSeek R1’s training methodology is Group Relative Policy Optimization (GRPO), a novel reinforcement learning technique specifically designed to enhance reasoning abilities [7]. GRPO optimizes the model’s capacity to develop coherent, logical reasoning patterns—a capability that distinguishes DeepSeek R1 in mathematical and complex problem-solving tasks [4][7][8].

This sophisticated, multi-phase training approach represents a significant departure from conventional training methodologies and has been instrumental in developing DeepSeek R1’s exceptional reasoning capabilities, particularly in domains requiring multi-step logical deduction [2][5][7].

Together, these four technical innovations—parameter efficiency through selective activation, the Mixture of Experts framework, advanced attention mechanisms, and breakthrough training methodologies—form the foundation of DeepSeek R1’s capabilities. They collectively represent not just incremental improvements to existing approaches but a fundamental rethinking of how large language models can be designed, trained, and deployed for maximum efficiency and performance.

### **III. Performance Profile: Benchmarks and Capabilities**

While DeepSeek R1’s technical architecture represents a significant innovation in model design, the ultimate measure of any AI system lies in its practical performance. In this section, we examine DeepSeek R1’s capabilities across four key domains, comparing its performance against leading models to assess its competitive positioning and disruptive potential in the large model industry.

#### **Mathematical Reasoning Excellence**

DeepSeek R1 demonstrates exceptional mathematical reasoning capabilities, outperforming even the most advanced models from industry leaders in several challenging benchmarks. On the American Invitational Mathematics Examination (AIME) 2024, DeepSeek R1 achieved a remarkable 79.8% accuracy, surpassing OpenAI’s o1-1217 model which scored 79.2% [4][7][8]. This performance is particularly noteworthy given that AIME represents one of the most challenging mathematical competitions for high school students in the United States, requiring sophisticated multi-step reasoning and problem-solving strategies.

Even more impressive is DeepSeek R1’s performance on the MATH-500 benchmark, where it achieved 97.3% accuracy compared to OpenAI o1’s 96.4% [4][7][8]. For context, GPT-4 achieves 98.2% on this benchmark [3], placing DeepSeek R1 in the highest tier of mathematical reasoning capabilities among current AI systems.

These mathematical reasoning capabilities stem directly from DeepSeek R1’s advanced training methodology, particularly its multi-phase approach that emphasizes chain-of-thought reasoning and the innovative Group Relative Policy Optimization technique [7]. The model’s ability to generate lengthy, detailed reasoning chains enables it to tackle complex mathematical problems that require multiple logical steps and sophisticated problem decomposition [2][5][7].

#### **Coding and Software Engineering Proficiency**

DeepSeek R1 demonstrates similarly impressive capabilities in coding and software engineering tasks. On the Codeforces benchmark, which evaluates competitive programming abilities,



DeepSeek R1 achieved 96.3% accuracy, nearly matching OpenAI o1-1217's 96.6% [4]. This performance indicates that DeepSeek R1 can generate sophisticated algorithms and solve complex programming challenges at a level comparable to the industry's leading models.

In the more practical domain of software engineering, DeepSeek R1 achieved 49.2% on the SWE-bench Verified benchmark, slightly outperforming OpenAI's o1-1217 at 48.9% [4]. SWE-bench evaluates a model's ability to understand, modify, and debug real-world codebases—tasks that more closely resemble the day-to-day work of software engineers than competitive programming challenges.

DeepSeek R1 also achieved 65.9% on the LiveCodeBench (Pass@1-COT) benchmark [5], further demonstrating its capabilities in generating functional code across multiple programming languages. These coding and software engineering capabilities position DeepSeek R1 as a powerful tool for developers, potentially accelerating software development processes and enabling more sophisticated code generation and analysis applications.

#### \*\*\* Multilingual Capabilities

While many large language models focus primarily on English, DeepSeek R1 demonstrates exceptional proficiency in both English and Chinese [1][7][9]. This bilingual capability represents a significant advantage in the global AI market, particularly for applications requiring cross-lingual understanding or deployment in Chinese-speaking regions.

DeepSeek R1's multilingual proficiency extends beyond simple translation to include nuanced understanding of cultural contexts, idioms, and domain-specific terminology in both languages. This capability stems from its diverse training data and sophisticated attention mechanisms that enable effective processing of different linguistic structures [10].

The model's strong performance across both English and Chinese benchmarks positions it uniquely in the global AI landscape, bridging Eastern and Western markets and potentially enabling more effective knowledge transfer and communication across linguistic boundaries.

#### \*\*\* Cost-Performance Analysis

Perhaps DeepSeek R1's most disruptive characteristic is its exceptional cost-performance ratio. With estimated training costs of approximately \$5.58 million—compared to GPT-4's reported \$100+ million—DeepSeek R1 achieves comparable performance at a fraction of the development investment [3][9]. This dramatic cost reduction stems directly from its parameter-efficient architecture, which requires significantly less computational resources during training while maintaining competitive performance.

Even more significant for potential users are DeepSeek R1's operational costs. At approximately \$8 per million tokens (for both input and output combined), DeepSeek R1 offers a substantial cost advantage over OpenAI's pricing structure, which can reach \$60 per million output tokens for their most advanced models [14]. This operational cost efficiency makes advanced AI capabilities accessible to a much wider range of organizations and use cases, potentially democratizing access to state-of-the-art AI technologies.

When considering the combination of performance metrics and cost efficiency, DeepSeek R1 emerges as a potentially disruptive force in the large model industry. It challenges the assumption that cutting-edge AI capabilities must come with prohibitive costs, instead demonstrating that innovative architectural approaches can deliver competitive performance at a fraction of the traditional resource requirements.

The model's exceptional mathematical reasoning, strong coding capabilities, multilingual proficiency, and revolutionary cost structure collectively position DeepSeek R1 as not just another incremental advancement in large language models, but potentially a paradigm shift in how these systems are designed, deployed, and commercialized. By achieving performance comparable to industry leaders at dramatically lower costs, DeepSeek R1 could accelerate AI adoption across industries and geographies, particularly in resource-constrained environments and emerging markets previously excluded from the benefits of advanced AI technologies.

## \*\* IV. Impact on the Large Model Industry

DeepSeek R1's emergence represents more than just another technical advancement in AI capabilities—it signals a potential paradigm shift in how large language models are developed, deployed, and commercialized. By challenging established norms around model architecture, pricing, and accessibility, DeepSeek R1 is reshaping standards, expectations, and practices across the AI industry. This section examines four key areas where DeepSeek R1's impact is particularly significant.

### \*\*\* Business Model Disruption

DeepSeek R1's revolutionary cost structure directly challenges the premium pricing models that have dominated the advanced AI market. With operational costs of approximately \$8 per million tokens (for both input and output combined), DeepSeek R1 undercuts OpenAI's pricing structure by up to 87%, where costs can reach \$60 per million output tokens for their most advanced models [14]. This dramatic price differential threatens to disrupt the revenue models of established players who have built their businesses around high-margin AI services [20].

The cost disruption extends beyond operational expenses to the initial investment required for model development. DeepSeek R1's training costs of approximately \$5.58 million—compared to GPT-4's reported \$100+ million—demonstrate that state-of-the-art AI capabilities no longer require massive capital expenditures [3][9]. This reduced barrier to entry enables new competitors to enter the market with lower initial investments, potentially accelerating innovation and competition in the AI space.

Perhaps most significantly, DeepSeek R1's open-source approach under the MIT License fundamentally challenges the proprietary business models that have dominated the industry [1][2][5][9]. This approach shifts the value proposition from exclusive access to foundation models toward specialized implementations, domain expertise, and complementary services. Organizations across the AI value chain are now forced to reconsider their competitive differentiation strategies, with many likely to focus on application-specific optimizations and industry-specific solutions rather than foundational model capabilities [25].

As the global LLM market expands from \$4.5 billion in 2023 to a projected \$82.1 billion by 2033 [19], DeepSeek R1's business model disruption could significantly reshape how this growth materializes, potentially accelerating adoption while compressing margins for basic AI services.

### \*\*\* Technical Standards Evolution

DeepSeek R1's parameter-efficient architecture is establishing new benchmarks for model evaluation and comparison. While total parameter count has historically been a primary metric for assessing model scale and capability, DeepSeek R1 shifts the focus toward the ratio of performance to activated parameters—a more meaningful measure of computational efficiency [3][7]. This new standard could fundamentally change how models are designed and evaluated, prioritizing selective computation over brute-force scaling [25].

The success of DeepSeek R1's Mixture of Experts (MoE) architecture is accelerating industry adoption of this approach. Major players are likely to incorporate similar techniques in their next-generation models, potentially establishing MoE as the standard architecture for large-scale AI systems [24]. This architectural shift addresses growing concerns about AI's environmental impact and operational costs while maintaining or improving performance capabilities [3][15].

DeepSeek R1's implementation of Multi-Head Latent Attention (MLA) with dynamic compression techniques is similarly influencing attention mechanism design across the industry [10][15]. As context length becomes increasingly important for complex reasoning tasks, these attention optimizations may become standard components in large model architectures, enabling more efficient processing of extended sequences without prohibitive computational requirements [1][2][5][11].

The hybrid training methodology employed by DeepSeek R1—combining supervised fine-tuning with multi-phase reinforcement learning—is reshaping industry standards for model training [2][4][5][7]. This methodical approach differs from more polarized methodologies currently

dominant and may establish a new benchmark for developing models with advanced reasoning capabilities [3][7][8].

#### \*\*\* Democratization of Advanced AI

DeepSeek R1's combination of cost efficiency and open-source availability is dramatically expanding access to advanced AI capabilities globally. Organizations previously priced out of the high-end AI market can now leverage state-of-the-art language models for applications ranging from content creation and data analysis to software development and customer service [9][17]. This democratization could accelerate AI adoption across industries and geographies, particularly in emerging markets and resource-constrained environments [14].

The model's tiered approach to deployment—offering distilled versions with varying parameter sizes—further enhances accessibility by accommodating diverse computational resources [2][4][5][7]. Organizations can select the appropriate model size based on their specific requirements and available infrastructure, from the 1.5B parameter distilled version suitable for edge devices to the full 671B parameter model for complex reasoning tasks [3][7].

This democratization aligns with broader industry trends toward more accessible AI technologies, as 67% of organizations are now utilizing generative AI products [19]. DeepSeek R1 could accelerate this adoption curve, particularly in sectors like retail, e-commerce, healthcare, and financial services, where task automation and operational efficiencies are key priorities [19][23].

#### \*\*\* Geopolitical Implications

As a Chinese AI startup achieving performance comparable to leading Western models, DeepSeek represents a significant development in the global AI landscape [20]. This success could accelerate the rebalancing of AI capabilities globally, potentially intensifying the technological competition between the United States and China—the two nations primarily driving the global AI race [20].

DeepSeek R1's emergence occurs against a backdrop of growing concerns about Chinese AI technologies in Western markets, with many organizations viewing Chinese AI companies as security risks due to data privacy, security, and censorship concerns [20]. These tensions could influence market access and regulatory responses, potentially fragmenting the global AI ecosystem along geopolitical lines.

The open-source nature of DeepSeek R1 adds another dimension to these geopolitical considerations. By making advanced AI capabilities widely accessible, DeepSeek challenges the ability of any single nation to control or restrict access to cutting-edge AI technologies [14]. This democratization could influence global AI governance frameworks, potentially accelerating international collaboration on AI safety, ethics, and regulation.

As North America maintains its dominant position with approximately 32-33% of the global LLM market share [19][23], and Asia Pacific—particularly China—makes significant AI investments [23], DeepSeek R1 represents a potential inflection point in this competitive dynamic. Its combination of performance, efficiency, and accessibility could reshape not just technical standards and business models, but the geopolitical balance of power in the global AI landscape.

#### \*\* VII. Future Development Directions

DeepSeek R1's innovative approach to large language model design signals several important directions for the future development of AI technologies. These emerging trends will likely shape how models are conceptualized, built, and deployed in the coming years.

#### \*\*\* Efficiency-Focused Innovation

DeepSeek R1's parameter-efficient architecture reinforces and accelerates the industry trend toward smaller, more efficient models that maintain high performance while reducing computational requirements [25]. This shift represents a fundamental rethinking of the "bigger is better" paradigm that has dominated AI development. Future models will likely prioritize architectural innovations that optimize the performance-to-resource ratio, potentially through more sophisticated routing mechanisms, enhanced expert specialization, and novel attention optimizations [15][25].

The success of DeepSeek R1’s selective activation approach—using only 37B of its 671B parameters per token—demonstrates that efficiency gains need not come at the expense of performance [3][7]. This realization could drive research toward even more selective computation strategies, potentially enabling models with trillions of total parameters but activating only a small fraction for any given task. Such developments would address growing concerns about AI’s environmental impact, operational costs, and accessibility [3][9].

#### \*\*\* Hardware-Software Co-Evolution

The success of DeepSeek R1’s MoE architecture will likely drive hardware innovations specifically optimized for these models [27][31]. Future AI accelerators may incorporate specialized components for efficient routing between experts, optimized memory hierarchies for selective parameter activation, and enhanced parallelization capabilities. This hardware-software co-evolution could further enhance the efficiency advantages of MoE models, potentially establishing them as the dominant architecture for large-scale AI systems [15].

Current hardware like NVIDIA’s H100, Intel’s Gaudi2, and AMD’s MI300X will likely evolve to better support the specific computational patterns of MoE architectures [27]. These developments might include enhanced memory bandwidth to support rapid expert switching, specialized tensor cores optimized for routing operations, and improved energy efficiency for sustainable AI deployment [31].

#### \*\*\* Collaborative Development Patterns

DeepSeek R1’s open-source approach under the MIT license could strengthen the trend toward collaborative AI development [1][2][5][9]. This model of innovation, where core capabilities are openly shared while companies compete on implementation, services, and specialized applications, may become increasingly prevalent [25]. Such collaboration could accelerate overall industry advancement while shifting competitive differentiation toward application expertise, domain knowledge, and service quality rather than foundational model capabilities.

#### \*\* VIII. Challenges and Considerations

Despite its promising capabilities, DeepSeek R1 faces several significant challenges that could impact its adoption and influence on the AI landscape.

##### \*\*\* Regulatory and Ethical Dimensions

The democratization of advanced AI capabilities through cost-effective, open-source models like DeepSeek R1 raises important questions about responsible use, misuse prevention, and appropriate governance frameworks [14]. As these powerful technologies become more widely accessible, regulatory bodies worldwide may implement stricter oversight mechanisms, potentially creating compliance challenges for organizations deploying these models [20].

Ethical considerations around bias, fairness, and transparency become increasingly important as AI deployment expands across diverse contexts and cultures. DeepSeek R1’s multilingual capabilities, while beneficial for global applications, may also introduce complex cross-cultural ethical considerations that require careful navigation [9].

##### \*\*\* Technical Limitations

Despite its impressive capabilities, DeepSeek R1 still faces technical limitations. Its performance on certain general knowledge benchmarks, such as GPQA Diamond (71.5% vs. OpenAI o1-1217’s 75.7%) and MMLU (90.8% vs. o1-1217’s 91.8%), indicates room for improvement in some domains [4][8]. Additionally, the complexity of the MoE architecture may present challenges for deployment in resource-constrained environments, potentially limiting its applicability in certain contexts despite the availability of distilled versions [3][7].

##### \*\*\* Adoption Barriers

As a Chinese AI technology, DeepSeek R1 may face adoption resistance in Western markets due to geopolitical tensions and security concerns [20]. Many Western organizations view Chinese AI companies as security risks due to data privacy, security, and censorship concerns, potentially limiting DeepSeek R1’s global market penetration despite its technical advantages [20].

Technical integration challenges also present potential barriers to adoption. Organizations with existing AI infrastructure optimized for traditional transformer architectures may face significant costs and complexities when transitioning to MoE-based models like DeepSeek R1, potentially slowing adoption despite the long-term efficiency benefits [26].

These challenges and considerations will play a crucial role in determining DeepSeek R1's ultimate impact on the AI landscape. How effectively these issues are addressed—by DeepSeek, the broader AI community, and regulatory bodies—will significantly influence whether DeepSeek R1 represents a temporary innovation or a lasting paradigm shift in large language model development.

#### \*\* IX. Conclusion: DeepSeek R1's Place in AI Evolution

DeepSeek R1 represents a potential inflection point in the evolution of artificial intelligence—one that challenges fundamental assumptions about how large language models should be designed, deployed, and commercialized. Throughout this analysis, we've seen how DeepSeek R1's innovative architecture marks a paradigm shift from the "bigger is better" approach that has dominated AI development toward a new focus on parameter efficiency and selective computation [3][7][25].

By activating only 37 billion of its 671 billion parameters per token, DeepSeek R1 demonstrates that cutting-edge AI performance need not come at the cost of computational efficiency [1][2][3]. This architectural breakthrough, combined with advanced attention mechanisms and sophisticated training methodologies, enables performance comparable to models with far greater resource requirements while dramatically reducing both development and operational costs [3][9][14].

Perhaps equally significant is DeepSeek's decision to release R1 under the MIT License, making advanced AI capabilities accessible to researchers, developers, and organizations worldwide [1][2][5][9]. This open-source approach, coupled with the model's revolutionary cost structure, has profound implications for the democratization of AI technologies—potentially accelerating adoption across industries and geographies previously excluded from the benefits of advanced AI [14][20].

Looking forward, DeepSeek R1 may herald a new phase in AI development characterized by efficiency-focused innovation, hardware-software co-evolution and collaborative development patterns [25][33]. As the industry increasingly recognizes the limitations of simply scaling up existing architectures, approaches that prioritize selective computation and domain-specific optimization are likely to gain prominence [21][25].

The ultimate impact of DeepSeek R1 will depend on how effectively it navigates technical limitations, regulatory considerations, and geopolitical tensions [20]. Yet regardless of its specific market trajectory, DeepSeek R1's technical innovations and disruptive approach have already challenged industry norms and expanded our understanding of what's possible in large language model design. In doing so, it may have permanently altered the course of AI evolution—shifting focus from how large models can become to how intelligently they can utilize their resources [3][25].