

Textual-to-Visual Iterative Self-Verification Slide Generation Agent

Anonymous ACL submission

Abstract

Generating presentation slides is a time-consuming task that urgently requires automation. Due to their limited flexibility and lack of automated refinement mechanisms, existing autonomous LLM-based agents face constraints in real-world applicability. In this work, we decompose the task of generating missing presentation slides into two key components: **content generation** and **layout generation**, aligning with the typical process of creating academic slides. For content generation, we introduce a content generation approach that enhances coherence and relevance by incorporating context from surrounding slides and leveraging section retrieval strategies. For layout generation, we propose a **textual-to-visual self-verification process** using a **LLM-based Reviewer + Refiner workflow**, transforming complex textual layouts into intuitive visual formats. This modality transformation simplifies the task, enabling accurate and human-like review and refinement. Experiments show that our approach significantly outperforms baseline methods in terms of alignment, logical flow, visual appeal, and readability.

1 Introduction

Effectively summarizing and presenting research findings through academic presentation slides is an essential part of scientific communication, allowing researchers to present key contributions and engage audiences at conferences and seminars (Guo et al., 2024; Mondal et al., 2024). However, creating these slides is a time-consuming process that requires extracting core information from lengthy papers, organizing it coherently, and designing visually consistent layouts across multiple slides (Fu et al., 2021). With the rapid growth in the volume of research, the demand for automated solutions has increased significantly. Recent advances in large language models (LLMs) (OpenAI, 2023; Touvron et al., 2023; Templeton et al., 2024) have

demonstrated remarkable capabilities in mimicking human behavior for complex tasks (Hong et al., 2023; Park et al., 2023; Yao et al., 2022b; Zala et al., 2024; Ma et al., 2024a) beyond text generation (Yao et al., 2022b,a; Xi et al., 2024; Yang et al., 2024). Building on these strengths, LLM-based agents offer a promising opportunity to automate tasks like slide generation (Zheng et al., 2025), reducing manual effort while ensuring coherence and visual quality.

Despite its potential, generating high-quality academic presentation slides presents two major challenges: **how to assign reasonable and adaptive layouts for generated content** and **how to ensure layout quality and consistency**.

The first challenge lies in generating layout information that adapts to the unique visual structure for different textual contents. Some methods focus solely on textual content, neglecting structural aspects like positioning, spacing, and alignment, leading to impractical outputs (Sun et al., 2021; Bandyopadhyay et al., 2024). Existing rule-based methods provide a quick and straightforward solution by populating predefined slots with generated content. However, they overlook the unique structural style of each presentation, often leading to rigid layouts that break the visual coherence.

The second challenge lies in achieving consistent textual-visual results, complicated by the inherent difficulty of representing slide layouts in structured textual formats. Unlike visual representations, where spatial relationships and element alignment are easy to interpret, textual formats lack this visual clarity (Xu et al., 2024; Hu et al., 2024). This makes it difficult for models to fully understand the spatial and structural aspects of slide design, leading to frequent errors such as text overflow, misalignment, and inconsistent spacing.

Furthermore, correcting these errors directly in the textual format is non-trivial. Without a visual reference, detecting overlapping elements or mis-

alignments becomes challenging, particularly in slides with complex layouts.

The key component of our framework is a textual-to-visual iterative self-verification process to refine initial output. The initial slide layouts are generated in a textual format, which—while structured and machine-readable—often contains errors due to the complexity of representing slide information in a non-visual form. Additionally, reviewing and refining these layouts in their original format is challenging and unintuitive. To address this, we introduce a **modality transformation** (Li et al., 2025) that converts the textual format into a visualized form. This transformation significantly reduces the complexity of the task, making it easier for the LLM-based Reviewer + Refiner workflow to detect and correct issues such as alignment and text overflow in a human-like, intuitive manner. The reviewer provides feedback by analyzing the visual representation of the slide layout. The feedback is then passed to the refiner, who applies the suggested adjustments to the structured layout in textual format. This iterative refinement process ensures higher-quality final outputs with improved coherence and visual consistency.

Our key contributions are as follows.

1. An agentic framework for slide generation including content and layout generation approaches, ensuring thematic consistency and visual coherence.
2. A textual-to-visual iterative self-verification process with modality transformation, enabling intuitive and accurate refinement for slide layout.
3. Extensive analyses and systematic evaluation, demonstrating the significant effectiveness and practical potential of our framework for automated academic slide generation.

2 Related Work

In this section, we introduce the background of the LLM-based agent and existed studies on slides generations.

2.1 LLM-based Agent

LLMs have demonstrated impressive capabilities for complicated, interactive tasks (Yao et al., 2022b,a; Xi et al., 2024; Yang et al., 2024; Ma et al., 2024b). LLM-based autonomous agents have achieved remarkable progress in a wide range of domains, including logic reasoning (Qi et al., 2024; Khattab et al., 2022), tool use (Qin et al., 2024;

Zhang et al., 2023a), and social activities (Park et al., 2023). The current paradigm of agents relies on the language intelligence of LLMs. The mainstream work pattern encompasses environment perceiving, planning, reasoning, and executing, forming a workflow to dive and conquer intricate challenges.

Empowered by the recent progress of multi-modal pre-training, those agents can understand image, video, and audio channels (Wu et al., 2023; Liu et al., 2023). (i) Visual knowledge can largely facilitate reasoning and is integrated into Chain-of-Thoughts (Zhang et al., 2023b; Xu et al., 2024). (ii) Multi-modal reasoning enables divergent thinking cross modalities and takes advantage of those different modalities. Sketchpad (Hu et al., 2024) allows LLMs to draw drafts to assist its planning and reasoning, i.e., to draw auxiliary lines for geometry problems. Visualization-of-Thought (Wu et al., 2024) generates visual rationales for spatial reasoning tasks like mazes. For each stage of complex multi-modal tasks, selecting an appropriate modality as the main modality for reasoning can leverage the natural characteristics of the modality and stimulate the potential of LLMs (Park et al., 2025).

2.2 Slide Generation

Previous studies have explored extractive methods and simplified this task as sentence selection, e.g., to calculate the importance score and extract top sentences (Wang et al., 2017). With the development of small language models (Lewis et al., 2020; Raffel et al., 2020), slide generation is unified as abstractive, query-based document summarization (Sun et al., 2021).

Despite their early success, the emergence of LLMs exhibits exceptional performance and stimulates the demands of intelligent slide generation. Slide generation poses intricate challenges for autonomous agents, as it requires document reading comprehension and precise tool use to generate layouts. Pioneer work focuses on modifying target elements, asking agents to execute a series of specific instructions (Guo et al., 2024). The agent needs to understand the status of the slide, navigate to the element, and generate precise API calls. Recent studies first plan the outlines and then generate each page. To further control the style of presentations, Mondal et al. (2024) introduce a reward model trained on human feedback to guide both topic generation and content extraction. Consid-

er the visual quality of slides, Bandyopadhyay et al. (2024) employ a visual LM to insert images. DOC2PPT (Fu et al., 2021) integrates an object placer to predict the position and size of each element by training small models. PPTAgent (Zheng et al., 2025) directly utilizes slide templates to fix the layout and then fill textboxes, ensuring visual harmony and aesthetic appeal.

3 Methodology

In this section, we propose an LLM-based agentic workflow to automate the generation of content and layout for academic paper slides.

3.1 Task Formulation

We first formally define our slide generation task. In this task, a presentation is represented as a collection of slide pages, where each page consists of multiple elements. Each element $e \in E$ is a tuple (c, l) , where c denotes the content (e.g., text, images, tables) and l specifies the corresponding layout information (e.g., position, size, font style).

Our **overall task** is to generate the missing slide \hat{S}_i given the paper D , the missing slide topic T , and the partially available slide set $S = \{S_1, S_2, \dots, S_n\}$.

Input The input consists of: 1. A paper $D = \{d_1, d_2, \dots, d_m\}$, where d_i denotes a section or paragraph in the paper. 2. A missing slide topic T , describing the main focus of the missing slide. 3. A partially available slide set $S = \{S_1, S_2, \dots, S_n\}$, where some slides \hat{S}_i are missing. 4. The preceding slide S_{prev} and the following slide S_{next} as contextual information.

Output The output is a structured textual file \hat{S}_i , which describes the missing slide, including both content c and layout information l for each element $e \in E$. Formally,

$$\hat{S}_i = \{e_j = (c_j, l_j) \mid j = 1, 2, \dots, k\}$$

where k is the number of elements in the generated slide. The generated textual file can be directly converted into a PowerPoint slide.

3.2 Slide Generation Framework

The process of creating a presentation typically involves two key stages: (1) identifying the core content that needs to be presented on each slide, and (2) arranging this information into a visually coherent and consistent layout.

The goal of content generation is to generate c_j for each element e_j based on the paper D , the missing slide’s topic T , and contextual information from the surrounding slides S_{prev} and S_{next} :

$$c_j = \mathcal{G}_{\text{content}}(D, T, S_{prev}, S_{next})$$

Here, $\mathcal{G}_{\text{content}}$ represents the content generation process, ensuring that the generated content is accurate, concise, and contextually relevant.

The layout generation task determines the layout l_j for each element $e_j = (c_j, l_j)$ to maintain visual consistency and readability. The initial layout draft $l_j^{(0)}$ is generated using the content c_j and contextual information from the surrounding slides:

$$l_j^{(0)} = \mathcal{G}_{\text{layout_draft}}(c_j, S_{prev}, S_{next})$$

To refine the initial layout, a textual-to-visual iterative self-verification process is applied. The layout at step k ($l_j^{(k)}$) is visualized as $\text{Image}(l_j^{(k)})$, allowing the LLM-based Reviewer + Refiner workflow to provide feedback and corrections:

$$l_j^{(k+1)} = \mathcal{G}_{\text{refine}}(l_j^{(k)}, \text{Image}(l_j^{(k)}))$$

This iterative process continues until the layout reaches the desired quality and visual coherence.

3.2.1 Content Generation

Determining the key contents on a slide page involves understanding paper structures, extracting critical texts and figures, and ensuring overall coherence for a logical flow and consistent style.

Our content generation stage adopts a multi-step process with three sub-modules: Text Retriever, Figure Extractor, and Content Generator, consisting of a pipeline to identify relevant text segments, recommend figures and tables, and then decide the contents to present.

Text Retriever We build a text retriever to retrieve the most relevant sections of the paper. The paper is divided into section-level granularity, with each segment represented and indexed as a dense embedding. Given the topic of a slide, the retriever selects the most relevant segments by calculating the cosine similarity between the dense embeddings of the slide topic and the indexed sections.

Figure Extractor Beyond the retrieved text, the Figure Extractor identifies candidate figures to support slide content. It scans the top- k retrieved text segments for explicit references (e.g., “Figure 1”,

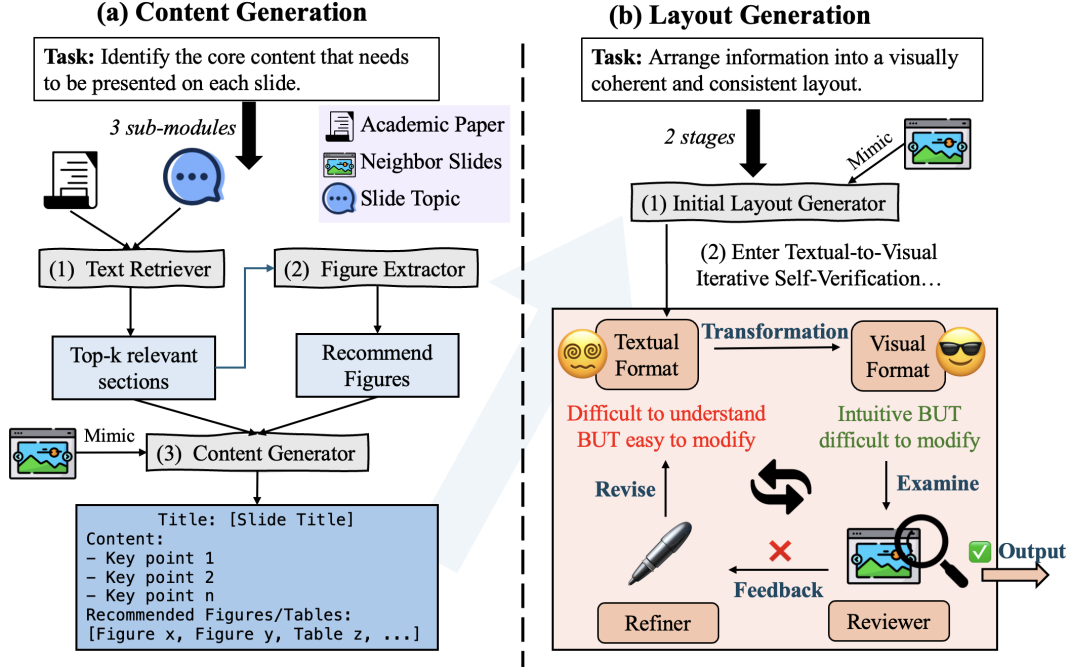


Figure 1: Overall Framework

“Table 2”), and extracts the corresponding captions from the paper. These captions provide semantic descriptions of the figures.

At this stage, only candidate figures are collected; the final selection is made by the Content Generator based on the generated slide content.

Content Generator Given the related text segments and candidate figures provided by previous modules, the LLM agent performs three sub-tasks.

First, it generates slide text aligned with the slide’s topic and context.

Second, it selects the most relevant figures or tables from the candidate set.

Finally, it incorporates adjacent slides to maintain logical flow and ensure seamless transitions.

The Content Generator’s output is then used for Layout Generation, which focuses on organizing content into a visually coherent slide.

3.2.2 Layout Generation

Slide layouts need to be flexible and controllable, rather than fully randomized or constrained by rigid templates. However, generating adaptive layouts is challenging and prone to issues such as text overflow, misalignment, and inconsistent spacing, especially when handling diverse content and styles.

To address this, we design a **textual-to-visual iterative self-verification process**. The initial layout draft mimics surrounding slides for style consistency but remains difficult to review in its struc-

tured textual format. We design an LLM-based *Reviewer* + *Refiner* workflow that validates and refines the layout respectively, improving accuracy and coherence through iterative corrections.

Stage 1: Initial Layout Generation The initial attempt is conducted by directly asking the LLM to arrange the layout for each element of the generated contents, specifying each element’s position, size, font, and color. We also append surrounding slide pages as demonstrations and carefully optimize the prompt to instruct the LLM to mimic their layout patterns for a visually consistent design. The layout is normalized as a JSON format.

While this initial layout serves as a foundation, our pilot experiments show that several factors contribute to potential errors:

(i) Textual slide layout is inherently complex, requiring detailed key-value pairs for positions, sizes, fonts, and colors. Any inconsistency in this structured data can cause significant visual defects.

(ii) LLMs lack direct visual feedback and cannot accurately assess how the generated layout will appear in its final form. Unlike models specifically trained for visual tasks, LLMs rely on textual context and structural patterns to predict layout information. This process is inherently limited, as it depends heavily on imitation and pattern recognition without understanding visual balance or spatial relationships. Consequently, the generated layouts

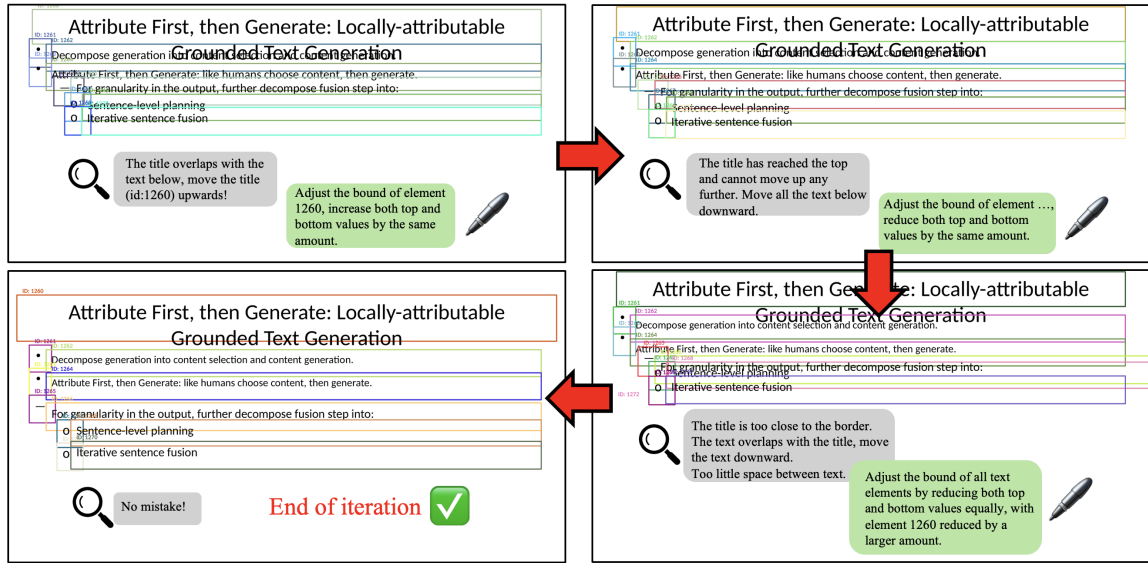


Figure 2: Iterative Layout Refinement in the Reviewer + Refiner Workflow

may exhibit issues such as poor alignment, overlapping elements, or inconsistent spacing, which require further refinement to ensure high-quality results.

Stage 2: Textual-to-Visual Iterative Self-Verification To refine the initial layout, we introduce a self-verification process that combines modality transformation and a LLM-based agentic workflow.

Modality Transformation We first convert the initial textual output into a visualized slide. The initialized layout is written into a slide and saved as an image. To facilitate visual perception, each visualized element in the slide is enclosed in a colored bounding box with a unique **ID**, matching its corresponding element in the textual file. This visual augmentation simplifies the workload, largely relieving the burden of perception and enabling the Reviewer to quickly reference specific elements and detect potential issues.

Reviewer The Reviewer simulates how a human expert would evaluate slide quality, following a predefined set of evaluation criteria and adjustment rules. Specifically, it performs the following tasks: Object overlapping detection, Image quality and distortion analysis, Element bounding and text overflow correction, Element positioning and alignment, Text formatting consistency and Overall composition and visual balance

Each recommendation is output as a structured list of suggestions, identifying specific elements

by their **ID** and providing precise numerical values for adjustments. For example, the Reviewer might suggest increasing a text box’s height by 1.2x to accommodate overflowing text or shifting an image downward by 10% of its height to resolve an overlap. Such a definite, specific advice format makes it easier for the Refiner to implement precise corrections in the subsequent refinement stage.

Refiner The Refiner plays a role for execution, translating the Reviewer’s visual feedback into precise modifications within the textual layout. To ensure accurate modifications, the Refiner follows a set of predefined rules based on the type of feedback received. For example, when the Reviewer suggests repositioning an element, the Refiner adjusts its bounding box coordinates accordingly while ensuring it remains within slide boundaries. Each rule is applied systematically based on the Reviewer’s feedback. The Refiner’s task is to modify only the necessary fields while maintaining the basic structure, resulting in a complete and refined file that reflects the intended adjustments.

Integration and Rendering The final output of this process is a refined JSON-formatted layout description that accurately represents the corrected slide. This JSON is passed to the rendering module to produce the final PowerPoint slide, ensuring that the layout visually reasonable and aligns with the overall presentation style.

LLM	Method	Coverage	ROUGE-1			ROUGE-2			ROUGE-L		
			P	R	F1	P	R	F1	P	R	F1
-	D2S	24.38	18.30	30.31	20.47	4.73	7.79	5.26	16.86	27.21	19.08
	Vanilla	29.81	24.56	47.74	28.02	8.94	19.96	10.34	17.54	37.58	20.46
	- w/o Retriever	28.18	30.06	42.04	29.35	12.44	19.45	12.54	23.19	34.85	22.99
	- w/o Neighbor	32.36	25.31	42.31	26.79	9.78	19.03	10.72	19.00	34.07	20.42
	Ours (3S)	32.76	28.64	39.30	27.47	11.23	17.13	11.15	21.99	32.18	21.36
	Ours (5S)	31.18	28.52	42.63	28.40	11.38	19.33	11.68	21.76	34.99	21.97
GPT-4o	Vanilla	28.41	23.29	43.97	25.65	7.15	16.86	8.20	16.23	34.09	18.31
	- w/o Retriever	29.01	32.48	37.68	28.36	11.15	15.88	10.05	24.45	30.35	21.64
	- w/o Neighbor	28.81	29.11	34.60	26.13	10.18	15.43	9.61	22.79	29.21	20.88
	Ours (3S)	29.49	31.63	32.86	26.10	11.30	14.91	9.84	24.34	27.81	20.76
	Ours (5S)	29.41	31.75	37.68	28.39	10.89	15.71	10.28	24.09	30.60	21.97
Qwen2.5-7B	Vanilla	25.28	24.27	44.92	26.02	9.06	19.69	10.10	17.89	36.24	19.65
	- w/o Retriever	26.18	31.47	36.77	27.92	12.60	17.11	11.60	24.66	30.39	22.14
	- w/o Neighbor	30.08	24.13	44.93	25.91	9.01	19.69	10.06	17.78	36.26	19.57
	Ours (3S)	28.79	29.78	36.26	25.99	11.63	16.58	10.56	24.17	30.76	21.21
	Ours (5S)	27.67	28.31	37.17	26.01	10.29	15.71	9.87	21.60	30.21	20.18

Table 1: Evaluation results for content generation.

4 Experiments

4.1 Dataset Construction

The dataset is sourced from the ACL 2024 In-Person Poster Session 1, with data collected from the public academic platform [Underline](#). The dataset consists of academic papers and their corresponding PowerPoint slides in PDF format, covering various research topics in natural language processing. To facilitate processing and preserve format details, all data is uniformly converted into JSON format, containing element-level information such as text content, font styles, positions, and sizes. Text from papers was extracted using GRO-BID ([Kermitt2, 2020](#)). Figures and captions were extracted using PDFFigures 2.0 ([Clark and Divvala, 2016](#)).

4.2 Baseline

We compare with two traditional document-to-slide generation systems.

D2S ([Sun et al., 2021](#)) adopts a two-step pipeline that first retrieves content using slide titles, then summarizes the retrieved content into bullet points via long-form QA.

Doc2ppt ([Fu et al., 2021](#)) formulates the task as end-to-end generation using a hierarchical seq2seq model that jointly predicts content and layout.

4.3 Implementation

We compare the performance of three large language models: **Llama-31-8B-Instruct** ([Grattafiori et al., 2024](#)), **GPT-4o** ([OpenAI et al., 2024](#)), and

Qwen-2.5-7B ([Qwen et al., 2025](#)). The best-performing model is selected to generate the final structured content. In the layout generation module, both the Reviewer and Refiner modules are built on top of multimodal large language model.

For the retriever, we use the **Salesforce SFR-Embedding-Mistral** ([Wang et al., 2024](#)) retriever to compute similarity scores and select the top-k most relevant sections.

Our experiments are naturally organized in the form of ablations. In the **w/o Section Retriever** configuration, the model receives the entire paper as input without section-level retrieval. In the **w/o Neighbor Slides** configuration, the surrounding slide content is removed, which helps assess the role of contextual information in maintaining logical flow and consistency.

4.4 Evaluation

Our evaluation method measures both content generation and layout generation. The evaluation process combines quantitative metrics and structured qualitative assessment to ensure comprehensive analysis.

Content Evaluation We evaluate generated slide content using ROUGE ([Lin, 2004](#)) and Coverage ([Kothawade et al., 2020](#)). ROUGE measures lexical overlap with reference slides, while Coverage computes the average cosine similarity between sentence embeddings from the source document and generated bullet points, reflecting semantic alignment. Higher Coverage indicates better semantic alignment with the original content.

Result Type	Element-Level	Slide-Level		Overall Impression	
	Align & Space	Logic	Coherence	Visual Appeal	Readability
Reference Slide	4.5	3.7	3.8	3.5	3.8
DOC2PPT	3.1	1.9	3.2	2.0	3.0
Baseline	2.0	3.0	3.3	2.0	2.5
JSON-Based Refinement	2.1	2.6	3.4	1.8	2.4
Our Method	3.0	3.8	3.4	2.8	3.1

Table 2: Evaluation results for layout generation

Layout Evaluation We adopt LLM-as-Judge (Chen et al., 2024) to evaluate slide layout across three levels:

- **Element Level:** Assesses alignment, spacing, and positioning of individual elements to ensure a well-structured layout.

- **Slide Level:** Focuses on logical flow and text-visual consistency, ensuring information is presented clearly and supported by relevant visuals.

- **Overall Impression:** Evaluates visual appeal and readability, ensuring cohesive design, appropriate font size, and clear figures.

4.5 Main Results

Content Generation Among the three models, GPT-4o demonstrates the most consistent and high performance, particularly in ROUGE-L F1 (21.97) and ROUGE-2 Recall (15.71). Although Llama-31-8B shows competitive performance in certain cases (e.g., ROUGE-1 Recall 47.74 for the Baseline), GPT-4o achieves a better balance between precision and recall. Qwen2.5-7B shows moderate performance, but its results are slightly more variable compared to the other models.

Layout Generation For layout evaluation, Table 2 summarizes the results of layout generation across three different configurations: Baseline, Textual-Based Refinement, and Our Method. The Reference Slide serves as a benchmark for assessing the quality of generated layouts.

Vanilla: This configuration represents the initial layout generated by the model without any refinement. The layout is stored in a structured JSON format describing positions, sizes, and other attributes. However, due to the complexity of multi-element layouts and the lack of visual feedback, this initial output often contains errors such as misalignment, overflow, and inconsistent spacing.

Textual-Based Refinement: In this configuration, the initial JSON file is refined through an automated rule-based review. The Reviewer analyzes

the JSON structure to detect layout issues, while the Refiner applies corrective actions directly to the JSON file. Although this approach improves some metrics, such as **Coherence (3.4)**, it still struggles with **Visual Appeal (1.8)** and **Alignment (2.1)**, indicating the limitations of rule-based refinement without visual feedback.

Our Method: By introducing **modality transformation**, we convert the JSON layout into a fully visualized slide image, allowing the Reviewer + Refiner workflow to detect and correct issues more intuitively. This approach yields significant improvements, especially in **Alignment and Spacing (3.0)** and **Logical Flow (3.8)**, closely approaching the quality of the reference slides. Additionally, **Visual Appeal (2.8)** and **Readability (3.0)** show notable gains compared to the previous configurations.

The results indicate that incorporating the Reviewer + Refiner workflow and modality transformation significantly improves layout quality, especially in terms of visual appeal and overall readability.

5 Analysis

5.1 Ablation

Effect of Neighbor Slides Neighbor slides significantly impact the quality of content generation. For instance, removing neighbor slides in Llama-31-8B (w/o Neighbor Slides) leads to a noticeable decrease in ROUGE-1 F1 (28.40 to 26.79) and ROUGE-2 F1 (11.68 to 10.72). Similar trends are observed in GPT-4o and Qwen2.5-7B, highlighting the importance of contextual information in maintaining logical coherence and reducing redundancy.

Balancing Full Context vs. Section Retrieval While using a section retriever helps reduce input length and improve efficiency, it can also cause minor variations in ROUGE scores. For example, Llama-31-8B with Section Retriever achieves slightly lower recall compared to its full-input coun-

terpart. When provided with the full paper, they can better understand the broader context and underlying relationships, resulting in more accurate and coherent slide content. This suggests that LLMs have strong capabilities in processing long documents. Thus, in scenarios where the input length remains within the allowable range, feeding the full paper is often more advantageous for generating high-quality slides on a given topic.

However, in situations where the input length exceeds the model’s context window or when the paper contains a significant amount of irrelevant information, **Section Retrieval** becomes essential. Selecting an optimal number of sections (e.g., 3 vs. 5) helps balance relevance and completeness. According to the results, **Ours (5S)** generally offers better recall and overall F1 compared to selecting fewer sections, as it provides more comprehensive contextual information without overwhelming the model with unnecessary details.

In summary, choosing between full-context input and section retrieval depends on the specific characteristics of the input paper. When the paper is relatively concise and highly relevant to the target topic, full-context input should be preferred. In contrast, for longer papers with diverse content, section retrieval is crucial for ensuring relevance while maintaining efficiency.

5.2 Factors Affecting Layout Quality

Alignment and Spacing metrics evaluate whether elements are properly positioned, evenly spaced, and free from overlap. As shown in Table 2, our method achieved a notable improvement in the Alignment and Spacing score (3.0) compared to the Baseline (2.0) and JSON-Based Refinement (2.1). Specifically, we observed that self-verification on textual layout cannot improve the layout quality, even compromise the Logic, Visual Appeal, and Readability. Our method eliminates this problem and achieves consistent improvement by introducing the textual-to-visual modality transformation.

Taking a closer look at the wrong cases, the remaining problems fall into three types.

(i) Low-quality initial layouts—such as overlapping elements or uneven spacing—limit the Reviewer’s ability to provide precise corrections. For example, when multiple elements overlap, it becomes unclear which one should be adjusted.

(ii) The lack of diverse layout patterns in the training data, particularly for slides with images, limits the model’s ability to position visual ele-

ments effectively.

(iii) Complex multi-element layouts can cause small errors to cascade during refinement, making them hard to fix without advanced optimization.

5.3 Complete Presentation Generation

While our current framework focuses on generating slides given a specific topic, the method can be naturally extended to automate the generation of a complete presentation composed of various slides.

Topic Generation and Slide Planning The first step in generating a full presentation is to extract key topics from the input paper. This can be achieved by analyzing the paper’s structure (e.g., Abstract, Introduction, Method, Results). Additionally, keyword extraction and clustering techniques can help create a sequence of logically connected topics for the slides. Each generated topic corresponds to a unique slide.

Multi-Page Content Generation Once the topics are generated, the framework applies the content generation strategy iteratively for each slide. By incorporating context from the previously generated slides, the model maintains logical flow and coherence across the entire presentation. Special transition slides (e.g., Overview) can be inserted to improve the presentation’s structure.

Consistent Layout and Visual Style The Reviewer + Refiner review process can be fully reused to ensure layout consistency across all slides.

This extension to full presentation generation holds significant practical value. It allows researchers to generate complete, high-quality presentations directly from academic papers, reducing the manual effort involved in slide creation.

6 Conclusion

In this paper, we propose a novel framework for generating academic presentation slides. By decomposing the task into content generation and layout generation, our method ensures adaptive layouts and visually consistent slides. We introduce a textual-to-visual iterative self-verification process using an LLM-based Reviewer + Refiner workflow, transforming complex textual layouts into visual representations for intuitive review and refinement. Experiments demonstrate that our approach significantly improves alignment, logical flow, visual appeal, and readability, offering a practical solution for automating high-quality slide generation.

633 Limitations

634 While our framework shows promising results in
635 generating academic slides, it has two main lim-
636 itations. First, the dataset is restricted to scien-
637 tific papers and corresponding presentation slides
638 from publicly available sources, which may limit
639 its generalizability to other types of presentations.
640 Second, the focus of our approach is primarily on
641 generating accurate content and structured layouts,
642 without considering advanced visual design aspects
643 such as color schemes, animations, or aesthetic en-
644 hancements that contribute to overall slide polish
645 and engagement.

646 Ethics Statement

647 Our study utilizes academic papers and their associ-
648 ated presentation slides that are publicly accessible
649 on official conference platforms. These materi-
650 als are collected under the fair use principle and
651 are used strictly for non-commercial, analytical
652 research purposes. We do not use these data for
653 model training, but only for performance evalua-
654 tion.

655 References

656 Sambaran Bandyopadhyay, Himanshu Maheshwari,
657 Anandhavelu Natarajan, and Apoorv Saxena. 2024.
658 [Enhancing presentation slide generation by LLMs
659 with a multi-staged end-to-end approach](#). In *Proceed-
660 ings of the 17th International Natural Language Gen-
661 eration Conference*, pages 222–229, Tokyo, Japan.
662 Association for Computational Linguistics.

663 Dongping Chen, Ruoxi Chen, Shilin Zhang, Yinyao Liu,
664 Yaochen Wang, Huichi Zhou, Qihui Zhang, Yao Wan,
665 Pan Zhou, and Lichao Sun. 2024. [Mllm-as-a-judge:
666 Assessing multimodal llm-as-a-judge with vision-
667 language benchmark](#). *Preprint*, arXiv:2402.04788.

668 Christopher Clark and Santosh Divvala. 2016. [Pdf-
669 figures 2.0: Mining figures from research papers](#).
670 In *Proceedings of the 16th ACM/IEEE-CS on Joint
671 Conference on Digital Libraries, JCDL '16*, page
672 143–152, New York, NY, USA. Association for Com-
673 puting Machinery.

674 Tsu-Jui Fu, William Yang Wang, Daniel J. McDuff, and
675 Yale Song. 2021. [Doc2ppt: Automatic presentation
676 slides generation from scientific documents](#). In *AAAI
677 Conference on Artificial Intelligence*.

678 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
679 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
680 Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-
681 ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh
682 Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-
683 tra, Archie Sravankumar, Artem Korenev, Arthur

Hinsvark, and 542 others. 2024. [The llama 3 herd of
models](#). *Preprint*, arXiv:2407.21783.

Yiduo Guo, Zekai Zhang, Yaobo Liang, Dongyan Zhao,
and Nan Duan. 2024. [PPTC benchmark: Evaluat-
ing large language models for PowerPoint task com-
pletion](#). In *Findings of the Association for Computa-
tional Linguistics: ACL 2024*, pages 8682–8701,
Bangkok, Thailand. Association for Computational
Linguistics.

Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng
Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang,
Yuxiao Dong, Ming Ding, and 1 others. 2023. [Coga-
gent: A visual language model for gui agents](#). *ArXiv
preprint*, abs/2312.08914.

Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Os-
tendorf, Luke Zettlemoyer, Noah A. Smith, and Ran-
jay Krishna. 2024. [Visual sketchpad: Sketching as
a visual chain of thought for multimodal language
models](#). In *The Thirty-eighth Annual Conference on
Neural Information Processing Systems*.

Kermitt2. 2020. Grobid: Machine learning for
extracting information from scholarly documents.
<https://github.com/kermitt2/grobid>. Ac-
cessed: 2025-02-16.

Omar Khattab, Keshav Santhanam, Xiang Lisa
Li, David Hall, Percy Liang, Christopher Potts,
and Matei Zaharia. 2022. Demonstrate-search-
predict: Composing retrieval and language mod-
els for knowledge-intensive nlp. *arXiv preprint
arXiv:2212.14024*.

Suraj Kothawade, Jiten Girdhar, Chandrashekhar Lava-
nia, and Rishabh Iyer. 2020. [Deep submodular net-
works for extractive data summarization](#). *Preprint*,
arXiv:2010.08593.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan
Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
Veselin Stoyanov, and Luke Zettlemoyer. 2020.
[BART: Denoising sequence-to-sequence pre-training
for natural language generation, translation, and com-
prehension](#). In *Proceedings of the 58th Annual Meet-
ing of the Association for Computational Linguistics*,
pages 7871–7880, Online. Association for Computa-
tional Linguistics.

Chengzu Li, Wenshan Wu, Huanyu Zhang, Yan Xia,
Shaoguang Mao, Li Dong, Ivan Vulić, and Furu
Wei. 2025. [Imagine while reasoning in space:
Multimodal visualization-of-thought](#). *Preprint*,
arXiv:2501.07542.

Chin-Yew Lin. 2004. [ROUGE: A package for auto-
matic evaluation of summaries](#). In *Text Summariza-
tion Branches Out*, pages 74–81, Barcelona, Spain.
Association for Computational Linguistics.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae
Lee. 2023. [Visual instruction tuning](#).

Xinbei Ma, Yiting Wang, Yao Yao, Tongxin Yuan, Aston Zhang, Zhuosheng Zhang, and Hai Zhao. 2024a. Caution for the environment: Multimodal agents are susceptible to environmental distractions. <i>arXiv preprint arXiv:2408.02544</i> .	795
Xinbei Ma, Zhuosheng Zhang, and Hai Zhao. 2024b. Coco-agent: A comprehensive cognitive mllm agent for smartphone gui automation. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 9097–9110.	796
Ishani Mondal, Shwetha S, Anandhavelu Natarajan, Aparna Garimella, Sambaran Bandyopadhyay, and Jordan Boyd-Graber. 2024. Presentations by the humans and for the humans: Harnessing LLMs for generating persona-aware slides from documents . In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2664–2684, St. Julian’s, Malta. Association for Computational Linguistics.	797
OpenAI. 2023. Gpt-4 technical report . <i>ArXiv preprint</i> , abs/2303.08774.	798
OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	799
Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simula- cra of human behavior. In <i>In the 36th Annual ACM Symposium on User Interface Software and Technology (UIST ’23)</i> , UIST ’23, New York, NY, USA. Association for Computing Machinery.	800
Simon Park, Abhishek Panigrahi, Yun Cheng, Dingli Yu, Anirudh Goyal, and Sanjeev Arora. 2025. Generalizing from simple to hard visual reasoning: Can we mitigate modality imbalance in vlms? <i>Preprint</i> , arXiv:2501.02669.	801
Zhenting Qi, Mingyuan Ma, Jiahang Xu, Li Lyna Zhang, Fan Yang, and Mao Yang. 2024. Mutual reasoning makes smaller llms stronger problem-solvers . <i>Preprint</i> , arXiv:2408.06195.	802
Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, dahai li, Zhiyuan Liu, and Maosong Sun. 2024. ToolLLM: Facilitating large language models to master 16000+ real-world APIs . In <i>The Twelfth International Conference on Learning Representations</i> .	803
Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan	804
Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, and 25 others. 2025. Qwen2.5 technical report . <i>Preprint</i> , arXiv:2412.15115.	805
Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	806
Edward Sun, Yufang Hou, Dakuo Wang, Yunfeng Zhang, and Nancy X. R. Wang. 2021. D2S: Document-to-slide generation via query-based text summarization . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1405–1418, Online. Association for Computational Linguistics.	807
Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, and 3 others. 2024. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet . <i>Transformer Circuits Thread</i> .	808
Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>ArXiv preprint</i> , abs/2307.09288.	809
Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. Improving text embeddings with large language models . <i>Preprint</i> , arXiv:2401.00368.	810
Sida Wang, Xiaojun Wan, and Shikang Du. 2017. Phrase-based presentation slides generation for academic papers . In <i>AAAI Conference on Artificial Intelligence</i> .	811
Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2023. Next-gpt: Any-to-any multimodal llm .	812
Wenshan Wu, Shaoguang Mao, Yadong Zhang, Yan Xia, Li Dong, Lei Cui, and Furu Wei. 2024. Mind’s eye of LLMs: Visualization-of-thought elicits spatial reasoning in large language models . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	813
Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, Songyang Gao, Lu Chen, Rui Zheng, Yicheng Zou, Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. 2024. Agentgym: Evolving large language model-based agents across diverse environments . <i>Preprint</i> , arXiv:2406.04151.	814

Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. 2024. [Llava-cot: Let vision language models reason step-by-step](#). *Preprint*, arXiv:2411.10440.

John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering. *arXiv preprint arXiv:2405.15793*.

Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022b. [ReAct: Synergizing reasoning and acting in language models](#). volume abs/2210.03629.

Abhay Zala, Han Lin, Jaemin Cho, and Mohit Bansal. 2024. Diagrammergpt: Generating open-domain, open-platform diagrams via llm planning. In *COLM*.

Zhuosheng Zhang, Yao Yao, Aston Zhang, Xiangru Tang, Xinbei Ma, Zhiwei He, Yiming Wang, Mark Gerstein, Rui Wang, Gongshen Liu, and Hai Zhao. 2023a. [Igniting language intelligence: The hitchhiker’s guide from chain-of-thought reasoning to language agents](#). *Preprint*, arXiv:2311.11797.

Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023b. Multi-modal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*.

Hao Zheng, Xinyan Guan, Hao Kong, Jia Zheng, Hongyu Lin, Yaojie Lu, Ben He, Xianpei Han, and Le Sun. 2025. Pptagent: Generating and evaluating presentations beyond text-to-slides. *arXiv preprint arXiv:2501.03936*.

A Detailed Descriptions of Reviewer and Refiner Modules

A.1 Reviewer Module

The Reviewer module analyzes the visual representation of the slide, identifies layout issues, and provides precise feedback for improvements. This feedback focuses on alignment, spacing, text overflow, and image distortion. The primary goal of the Reviewer is to detect errors and ensure that all elements are properly positioned and formatted for a visually coherent slide.

A.1.1 Evaluation Criteria and Feedback Rules

The Reviewer module evaluates slides based on the following criteria:

- **Object Overlapping:** Identifies overlapping elements and suggests repositioning or resizing to maintain separation.
- **Image Quality and Distortion:** Detects blurry or distorted images and recommends proportional scaling.
- **Element Bounding and Text Overflow:** Ensures text fits within its bounding box and suggests expanding the box or reducing font size.
- **Element Positioning and Alignment:** Checks alignment and spacing, adjusting misaligned elements to the nearest grid line.
- **Text Formatting Consistency:** Verifies font family and text hierarchy, ensuring the title is larger than body text.
- **Overall Composition and Visual Balance:** Evaluates symmetry and visual balance, recommending layout adjustments for better harmony.

A.1.2 Example Output

The output of the Reviewer module is a structured JSON list, detailing necessary modifications for each slide element.

```
[
  {
    "element": 302,
    "recommendation": "Increase text box height by 1.2x to fit overflowing text."
  },
  {
    "element": 303,
    "recommendation": "Move downward by 10% of its height to resolve overlap with ID 302."
  },
  {
    "element": 304,
    "recommendation": "Reduce font size by 2pt to fit within the bounding box."
  }
]
```

A.2 Refiner Module

The Refiner module applies the Reviewer’s feedback by modifying the structured layout described in JSON format. This module focuses on correcting bounding box positions, resizing elements, and preventing overlaps.

953	A.2.1 Input to the Refiner	B Layout Evaluation Criteria and	1001
954	The input to the Refiner module consists of the	Scoring Standards	1002
955	following components:		
956	<ul style="list-style-type: none"> • JSON File: Describes the position, size, font, 	This section provides a detailed explanation of the	1003
957	and content of each element on the slide.	evaluation criteria used to assess the quality of the	1004
958	<ul style="list-style-type: none"> • Reviewer’s Feedback: Provides detailed rec- 	generated slides. The evaluation process covers	1005
959	ommendations for modifying elements (e.g.,	multiple aspects of slide design, including align-	1006
960	move, resize, align).	ment, logical flow, text-visual consistency, visual	1007
961	<ul style="list-style-type: none"> • Slide Dimensions: Ensures all adjustments 	appeal, and readability. Each criterion is scored on	1008
962	remain within the boundaries of the slide.	a five-point scale from 1 (Poor) to 5 (Excellent).	1009
963	A.2.2 Modification Instructions		1010
964	The Refiner applies modifications based on the Re-	B.1 Alignment and Spacing	1011
965	viewer’s feedback, following these rules:	This criterion evaluates whether elements on the	1012
966	<ul style="list-style-type: none"> • Move an Element: Adjust the element’s 	slide are properly positioned, evenly spaced, and	1013
967	bounding box values to reposition it. Mod-	free from overlap. It ensures that the layout main-	1014
968	ify the top, bottom, left, and right values as	tains visual balance and clarity.	1015
969	required.	<ul style="list-style-type: none"> • 1 Point (Poor): Severe misalignment; text 	1016
970	<ul style="list-style-type: none"> • Resize or Scale an Element: Modify the 	overlaps with visuals, creating a chaotic lay-	1017
971	width and height of an element proportionally	out.	1018
972	while preserving its aspect ratio.	<ul style="list-style-type: none"> • 3 Points (Average): Most elements are 	1019
973	<ul style="list-style-type: none"> • Avoid Overlap: Ensure no two elements over- 	aligned, but minor misplacements exist.	1020
974	lap by repositioning or resizing conflicting	<ul style="list-style-type: none"> • 5 Points (Excellent): Perfect alignment and 	1021
975	elements.	spacing with a professional layout.	1022
976	<ul style="list-style-type: none"> • Maintain Slide Boundaries: Prevent ele- 	Example Output:	1023
977	ments from exceeding the slide’s width or	{	1024
978	height.	"reason": "Most elements are well-	1025
979	A.2.3 Example Input and Output	aligned, but the spacing between	1026
980	The following example illustrates how the Refiner	the title and body text is	1027
981	module processes input and produces a refined lay-	inconsistent.",	1028
982	out.	"score": 4	1029
983	{	}	1030
984	"element": 302,		1031
985	"Bounds": [100, 200, 300, 400],	B.2 Logical Flow	1032
986	"Font": {"size": 16},	This criterion assesses the logical sequence of con-	1033
987	"Text": "Sample Text"	tent, ensuring that the information presented in the	1034
988	}	slide is clear and structured for easy audience un-	1035
989		derstanding.	1036
990	{	<ul style="list-style-type: none"> • 1 Point (Poor): Disorganized content; key 	1037
991	"element": 302,	points do not follow a logical sequence.	1038
992	"Bounds": [100, 220, 300, 420],	<ul style="list-style-type: none"> • 3 Points (Average): Basic logical structure; 	1039
993	"Font": {"size": 14},	minor reordering could improve the flow.	1040
994	"Text": "Sample Text"	<ul style="list-style-type: none"> • 5 Points (Excellent): Seamless logical se- 	1041
995	}	quence with clear and structured information.	1042
996	By applying these refinements iteratively, the	Example Output:	1043
997	Refiner ensures that the final slide layout meets		
998	high visual and structural standards, resulting in an		
999	accurate and human-like output.		
1000			

1044	{	1090	B.5 Readability
1045	"reason": "The information is	1091	This criterion evaluates the readability and clarity
1046	structured logically, but the	1092	of the text and graphical elements, ensuring that all
1047	second point would be clearer if	1093	content is easily understandable.
1048	placed before the third.",		
1049	"score": 4		
1050	}		
1051			
1052	B.3 Text-Visual Consistency		
1053	This criterion evaluates the consistency between		
1054	text and visual elements such as images and charts.		
1055	It ensures that visuals effectively support the textual		
1056	information.		
1057	• 1 Point (Poor): Visuals are irrelevant or con-	1094	• 1 Point (Poor): Text is too small or has low
1058	tradict the text.	1095	contrast, making it unreadable.
1059	• 3 Points (Average): Somewhat aligned, but	1096	• 3 Points (Average): Generally clear, but some
1060	better integration is needed.	1097	areas need better contrast or spacing.
1061	• 5 Points (Excellent): Perfectly integrated vi-	1098	• 5 Points (Excellent): Highly readable with
1062	suals that reinforce the message.	1099	optimal font size, spacing, and contrast.
1063	Example Output:	1100	Example Output:
1064	{	1101	{
1065	"reason": "The visuals effectively	1102	"reason": "The text is clear, well-
1066	support the content, but the chart	1103	spaced, and maintains good
1067	could be labeled more clearly.",	1104	contrast. The charts are easy to
1068	"score": 4	1105	read and properly scaled.",
1069	}	1106	"score": 5
1070		1107	}
1071	B.4 Visual Appeal	1108	These evaluation criteria ensure a comprehen-
1072	This criterion assesses the overall aesthetic quality	1109	sive and structured assessment of the generated
1073	of the slide, focusing on color harmony, typography,	1110	slides. By adhering to these standards, the evalua-
1074	and visual balance.	1111	tion process becomes interpretable, consistent, and
1075	• 1 Point (Poor): Inconsistent styling; visually	1112	reliable.
1076	unappealing design.	1113	
1077	• 3 Points (Average): Basic but functional color		
1078	scheme; lacks enhancements.		
1079	• 5 Points (Excellent): Cohesive and visually		
1080	appealing design with engaging elements.		
1081	Example Output:	1114	C Reliability Verification of Layout
1082	{	1115	Generation Evaluation
1083	"reason": "The color scheme is	1116	To verify the reliability of our layout evaluation
1084	visually appealing and harmonious,	1117	framework (LLM-as-Judge), we conducted a hu-
1085	but the background contrasts too	1118	man evaluation on a randomly selected subset of
1086	strongly with the text.",	1119	generated slides. To ensure consistency with the
1087	"score": 4	1120	LLM-based evaluations, we provided human raters
1088	}	1121	with the same scoring rubric and descriptions used
1089		1122	by the LLM, including detailed explanations of
		1123	each criterion. Each dimension was rated on a 1–5
		1124	Likert scale, where 1 indicates poor performance
		1125	and 5 indicates excellent performance. We then
		1126	computed Pearson correlation coefficients between
		1127	human scores and the LLM-based assessments.
		1128	As illustrated in Figure 6, the average correla-
		1129	tion reached 0.6984, suggesting a strong agreement
		1130	between human judgments and LLM evaluations.
		1131	In particular, dimensions such as <i>visual appeal</i> and
		1132	<i>readability</i> achieved the highest consistency, with
		1133	correlations of 0.89 and 0.77, respectively. These
		1134	results support the use of LLM-as-Judge as a reli-
		1135	able proxy for human evaluation in layout quality
		1136	evaluation.
		1137	

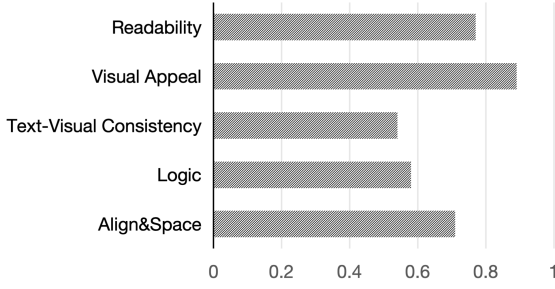


Figure 3: Pearson Correlation Between LLM-as-Judge and Human Evaluation Scores

D Comparison of Text Retrieval Strategy

To assess the impact of retrieval strategy, we compared the embedding-based retriever (Salesforce SFR Embedding-Mistral) with a classical sparse retrieval method, BM25. Both retrievers were applied to retrieve top-5 relevant segments for slide content generation. Quantitatively, the two approaches yielded comparable performance, and further analysis revealed a high degree of overlap in the retrieved segments. This is likely due to the formal and structured nature of academic writing, where key sentences often share significant lexical overlap—making sparse methods like BM25 surprisingly competitive.

Since final content generation is performed by a powerful LLM, minor differences in retrieval results tend to have limited influence on the final ROUGE scores. We retain the embedding-based retriever in our framework for its stronger generalization ability across domains and robustness in semantically diverse settings.

Retriever	ROUGE-1	ROUGE-2	ROUGE-L
BM25	28.67	11.27	22.49
Embedding-based	28.40	11.68	21.97

Table 3: F1 comparison of sparse and dense retrievers using LLaMA3-8B.

E Further Analysis by Slide Type

To better understand how our system performs across different types of slides, we conduct a qualitative analysis based on slide content composition. Specifically, we categorize generated slides into the following three types:

- **Text-only slides:** slides that contain only textual bullet points without any figures or tables.

- **Text + figure slides:** slides that combine textual content with at least one accompanying figure or table.
- **Figure-only slides:** slides where the primary content consists of visual elements, with minimal or no textual explanation.

For text-only slide, as shown in Figure 4, the refined version improves spacing between elements, making the content more legible and visually organized.

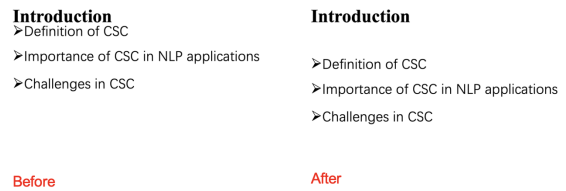


Figure 4: Example of Layout Refinement for a Text + Figure Slide

For text + figure slide, as shown in Figure 5, in the original version (left), the figure is relatively small and placed in the bottom-right corner, making it visually disconnected from the textual content. In the refined version (right), the figure is enlarged and repositioned to occupy the right half of the slide. The spacing between elements is also improved.

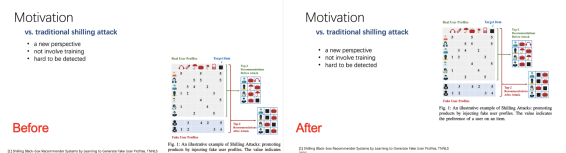


Figure 5: Example of Layout Refinement for a Text-only Slide

For figure-only slide, as shown in Figure 6, in the refined version (right), spacing and alignment are improved to reduce clutter and enhance readability, allowing each visual component to stand out more clearly.

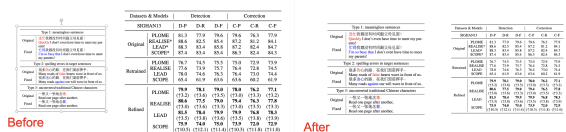


Figure 6: Example of Layout Refinement for a Figure-only Slide