

# Visual Executability: Enabling Language Models to Produce Executable Cinematic Plan

Anonymous ACL submission

## Abstract

Large language models (LLMs) have demonstrated strong capabilities in narrative understanding and generation, enabling applications such as script drafting and story planning. With the rapid rise of AI-generated web dramas and short-form narrative videos, there is increasing demand for systems that can translate narrative text into precise and controllable visual instructions at scale. However, LLM outputs are typically free-form and unstructured, making them difficult to integrate into downstream visual generation pipelines. While recent multimodal models attempt to directly map text to images or videos, such end-to-end approaches often lack interpretability, controllability, and compatibility with practical cinematic workflows.

In this work, we introduce *visual executability*, a property of language model outputs that enables deterministic execution by downstream visual systems. We propose a schema-guided framework that converts narrative text into shot-level executable cinematic plans, explicitly encoding scene context, character presence, emotional states, and cinematographic parameters such as shot type, camera movement, composition, and perspective. By constraining language model generation with a strict JSON schema, our approach guarantees structural validity, atomicity, and completeness of each visual unit, without requiring access to visual data or retraining multimodal models.

Through quantitative and qualitative evaluations on narrative texts, we demonstrate that our method produces fine-grained, coherent, and structurally valid cinematic plans that are directly usable for storyboarding and automated image or video generation pipelines. Our results suggest that enforcing visual executability at the language level provides a principled and modular alternative to end-to-end multimodal generation for controllable visual content creation.

## 1 Introduction

Large language models (LLMs) have demonstrated remarkable capabilities in narrative understanding and text generation, enabling applications such as creative writing, story summarization, and script drafting. Recent studies further show that LLMs can perform high-level reasoning and planning over complex semantic structures (Valmeekam et al., 2023; Guan et al., 2023; Pallagani et al., 2024). In parallel, AI-generated web dramas and short-form narrative videos have rapidly gained popularity, leading to the emergence of large-scale, fast-iteration content production pipelines. In these settings, narrative scripts must be translated into visual content efficiently and consistently, placing increasing pressure on automated systems to produce precise and controllable visual plans. Despite these advances, LLM outputs remain fundamentally non-executable with respect to visual generation systems. Bridging the gap between narrative language and controllable visual realization therefore remains a central challenge, particularly for applications such as storyboarding, film pre-visualization, and automated video generation.

Most existing approaches address this gap in one of two ways. Prompt-based methods rely on carefully engineered natural language prompts to guide text-to-image or text-to-video models. While flexible, these approaches produce free-form descriptions that lack explicit structure, making them difficult to verify, reuse, or systematically control (Zhou et al., 2023b). Alternatively, end-to-end multimodal models attempt to directly map narrative text to visual outputs through joint training on large-scale datasets. Although effective in certain scenarios, such models operate largely as black boxes and provide limited interpretability or explicit control over cinematographic elements such as shot boundaries, camera motion, and composition (He et al., 2023; Zhou et al., 2024; Yang et al.,

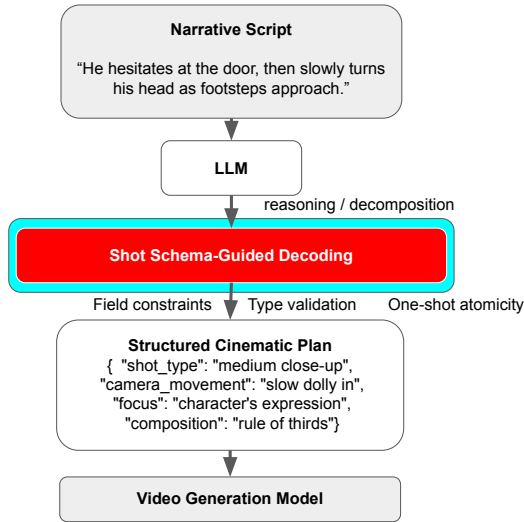


Figure 1: Overview of visual executability. Given a narrative script, a language model generates a structured cinematic plan through schema-guided decoding, enforcing field constraints, type validation, and one-shot atomicity. The resulting plan is directly executable by downstream video generation models.

2024).

In contrast, professional cinematic workflows do not generate visual content holistically. Instead, narratives are planned and executed at the level of shots—atomic visual units that specify what is shown, how it is framed, and how the camera moves. This shot-based planning paradigm underlies storyboards, shot lists, and pre-visualization pipelines, and relies on a shared visual language used by directors, cinematographers, and editors. Although recent work has explored language-based planning and action generation (Wang et al., 2023c; Zhou et al., 2023a), current language models do not natively produce outputs aligned with this shot-level visual representation.

Figure 1 summarizes the overall framework of our approach. Rather than generating free-form prompts or directly synthesizing visual outputs, we introduce an explicit intermediate representation in the form of a structured cinematic plan. Given a narrative script, a language model performs semantic reasoning and decomposition, while its generation process is constrained through schema-guided decoding. The schema enforces field completeness, type validity, and one-shot atomicity, resulting in a sequence of shot-level specifications that can be deterministically consumed by downstream image or video generation models.

Based on this formulation, we argue that the

key limitation of existing methods lies not in language understanding itself, but in the absence of an explicit intermediate representation that bridges narrative semantics and executable visual control. To address this limitation, we introduce the concept of *visual executability*, which characterizes whether a language model’s output can be deterministically mapped to a sequence of visual operations without additional interpretation or heuristic post-processing. Our notion of visual executability is closely related to recent advances in structured generation and constrained decoding (Wang et al., 2023a; Geng et al.; Zheng et al., 2024), but is specifically grounded in cinematic planning and visual realization.

We propose a schema-guided framework that transforms narrative text into executable cinematic plans composed of shot-level representations. Each shot is modeled as an atomic visual unit with explicitly defined scene context, character information, emotional tone, and cinematographic parameters such as shot type, composition, perspective, and camera movement. By enforcing a formal schema during generation, we guarantee structural validity by construction, eliminating the ambiguity inherent in free-form language.

Importantly, our approach does not require visual supervision or retraining of multimodal generation models. Instead, it leverages the implicit cinematic knowledge encoded in large-scale narrative corpora and makes this knowledge operational through structured generation. In this way, we shift the role of language models from passive narrators to active visual planners, capable of producing outputs that are both semantically grounded and operationally executable.

We evaluate our framework on narrative-to-cinematic planning tasks and demonstrate that it consistently produces fine-grained, coherent, and structurally valid cinematic plans. Human evaluations conducted with professional AI short-drama creators and film-directing students further indicate that the generated plans are directly usable for storyboarding and integration into automated image and video generation pipelines. Overall, our results suggest that visual executability provides a principled and modular alternative to end-to-end multimodal generation, enabling controllable and interpretable visual content creation driven by language models.

164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213

## 2 Related work

### 2.1 Language Models for Planning and Executable Representations

Recent studies have investigated the use of large language models (LLMs) for planning and decision-making, focusing on their ability to generate intermediate representations that support executable actions. A number of works analyze the planning capabilities of LLMs, revealing both their potential and inherent limitations in symbolic reasoning and long-horizon planning (Valmeekam et al., 2023; Pallagani et al., 2024). Other approaches leverage LLMs to construct explicit world models or planning domains that can be integrated into classical or model-based planning frameworks (Guan et al., 2023; Oswald et al., 2024).

Beyond symbolic planning, LLMs have been embedded into interactive or agent-based systems that unify reasoning, planning, and action selection in open-world environments (Wang et al., 2023c; Zhou et al., 2023a). Closely related work in embodied and language-to-action planning grounds language outputs in executable physical actions (Raman et al., 2022; Wu et al., 2023; Gramopadhye and Szafir, 2023). These methods define executability with respect to environment dynamics and task completion.

In contrast, our work focuses on visual executability rather than physical or symbolic action execution. We treat cinematic constructs—such as shots, camera movements, and framing—as executable units, and frame narrative understanding as a planning problem whose output is an intermediate representation designed specifically for visual realization.

### 2.2 Structured and Constrained Language Generation

Structured generation has emerged as a key technique for improving the reliability and controllability of language model outputs. Recent work evaluates schema- and grammar-constrained decoding, showing that formal constraints such as JSON schemas can significantly improve structural validity while maintaining output quality (Geng et al.; Wang et al., 2023a). Related approaches apply schema-guided generation to information extraction tasks, including scientific text and non-standardized tables, demonstrating that strict structural constraints enable robust downstream processing (Dagdelen et al., 2024; Hu et al., 2025).

More broadly, constrained and controlled text generation methods aim to steer LLM outputs using explicit instructions or structural requirements (Zhou et al., 2023b; Meng et al., 2022). Systems such as SGLang further emphasize the execution of structured language programs, highlighting the benefits of treating LLM outputs as executable artifacts rather than free-form text (Zheng et al., 2024). Feature Execution Graphs similarly formalize LLM outputs as executable graphs to support human–AI co-programming (Batatia and Svetlichnyi, 2025).

Our work builds on this line of research but departs in its objective. Rather than enforcing structure for semantic correctness or program execution, we impose schema constraints to guarantee visual executability, ensuring that each generated unit corresponds to an atomic and complete cinematic operation.

### 2.3 Visual Planning and Multimodal Generation from Language

Language-conditioned visual generation has been widely studied in both image and video domains. End-to-end story-to-image and story-to-video systems aim to directly synthesize visual content from narratives, often using diffusion-based or retrieval-augmented frameworks (He et al., 2023; Zhou et al., 2024; Yang et al., 2024). While these methods demonstrate impressive visual fidelity, they typically lack explicit intermediate representations for shot structure or cinematographic intent, limiting controllability and editability.

More recent work explores compositional and modular approaches to visual generation. Layout-GPT shows that LLMs can generate structured visual plans, such as layouts, that guide downstream image synthesis (Feng et al., 2023). Other studies emphasize the role of intermediate representations in multimodal systems, suggesting that structured control signals can simplify alignment and improve generation quality (Chen et al., 2025; Wang et al., 2023b).

Our approach aligns with these efforts but targets a different level of abstraction. Rather than generating spatial layouts or visual outputs directly, we propose a shot-level cinematic plan as an intermediate representation that captures temporal structure, visual focus, and cinematographic operations. This positions our work as a bridge between narrative language understanding and controllable multimodal generation.

Prior work has explored planning with language

214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264

models, structured generation, and multimodal synthesis from complementary perspectives. Our work differs by introducing visual executability as a unifying principle and by framing cinematic planning as a schema-constrained language generation problem. This enables language models to produce intermediate representations that are both semantically grounded and directly executable by visual generation pipelines.

## 3 Method

### 3.1 Problem Formulation

We study the problem of converting narrative text into a representation that can be deterministically executed by downstream visual generation systems. Let  $x \in \mathcal{X}$  denote an input narrative text. Our goal is to learn a mapping

$$f : \mathcal{X} \rightarrow \mathcal{P}, \quad (1)$$

where

$$\mathcal{P} = (s_1, s_2, \dots, s_N) \quad (2)$$

is an ordered sequence of cinematic shots. Each shot  $s_i$  represents an atomic visual unit that can be independently executed by a visual generation pipeline.

### 3.2 Visual Executability

We introduce *visual executability* as a property of language model outputs.

**Definition (Visual Executability).** A representation  $\mathcal{P}$  is said to be *visually executable* if there exists a deterministic mapping

$$g : \mathcal{P} \rightarrow \mathcal{A}, \quad (3)$$

where  $\mathcal{A}$  denotes a sequence of executable visual operations, such that no additional heuristic interpretation is required.

In our setting, visual executability is operationalized via the following criteria:

- **Atomicity:** each  $s_i$  corresponds to a single visual focus (e.g., one action, one reaction, or one observation).
- **Completeness:** each  $s_i$  explicitly specifies all attributes required for execution.
- **Discreteness:** visual attributes are drawn from predefined finite sets of primitives.
- **Structural Validity:**  $\mathcal{P}$  conforms to a formally defined schema that can be automatically validated.

### 3.3 Shot-Level Cinematic Representation

We model a cinematic plan as a sequence of shots, where each shot is treated as an *atomic visual action*. Specifically, each shot  $s_i$  is represented as a structured tuple:

$$s_i = \langle C_i, V_i, M_i, D_i \rangle, \quad (4)$$

where

$$\begin{aligned} C_i \in \mathcal{C} & \quad (\text{scene context}), \\ V_i \in \mathcal{V} & \quad (\text{visual focus and affect}), \\ M_i \in \mathcal{M} & \quad (\text{cinematographic parameters}), \\ D_i \in \mathcal{D} & \quad (\text{dialogue}). \end{aligned} \quad (5)$$

Intuitively,  $C_i$  captures contextual factors such as time and environment;  $V_i$  captures the visual focus (e.g., characters present) and the intended affect;  $M_i$  encodes cinematographic primitives such as shot type, camera perspective, composition, and camera movement; and  $D_i$  stores dialogue when present. This shot-level abstraction matches practical cinematic workflows and supports fine-grained decomposition of narratives into executable units.

We define a shot-level schema that specifies the required structure and attributes of each cinematic shot. Formally, each shot  $s_i$  must include the following fields:

- **Scene:** a textual description of the scene context (e.g., time, location, environment).
- **Characters:** a list of entities appearing in the shot.
- **Mood:** a categorical variable describing the emotional tone of the shot.
- **Shot Type:** a discrete cinematographic primitive (e.g., close-up, medium shot, long shot).
- **Composition:** a description of spatial composition (e.g., rule of thirds, centered).
- **Perspective:** the camera viewpoint (e.g., eye-level, low-angle, high-angle).
- **Camera Movement:** a discrete camera operation (e.g., static, pan, tilt, dolly).
- **Focus:** the primary visual focus of the shot.
- **Dialogue:** spoken content, if present.

All fields except *Dialogue* are required. Each field is associated with a predefined type and, where applicable, a finite set of admissible values.

### 3.4 Schema-Guided Generation

To guarantee visual executability, we constrain generation to a formally defined schema. Let  $\mathcal{S}$  denote a predefined schema specifying required fields, types, and admissible values for each shot and for the overall plan. We generate shots under schema constraints:

$$s_i \sim p_\theta(s | x, \mathcal{S}), \quad (6)$$

where  $p_\theta$  is restricted to outputs that satisfy  $\mathcal{S}$ . Equivalently, we enforce the hard constraint

$$\forall s_i \in \mathcal{P}, \quad s_i \in \mathcal{S}, \quad (7)$$

which ensures structural validity by construction and eliminates structurally invalid generations.

Compared to unconstrained free-form generation, schema-guided decoding shifts the model from descriptive narration to structured cinematic planning, producing outputs that can be deterministically mapped to downstream visual operations.

### 3.5 Implementation Note

In practice, the schema can be instantiated as a strict JSON schema with enumerated cinematographic primitives and required fields. This enables automatic validation of  $\mathcal{P}$  and supports modular integration with downstream storyboard tools and automated image/video generation pipelines.

## 4 Experiments

### 4.1 Experimental Setup

**Tasks and Data.** We evaluate cinematic planning on story-to-video planning tasks, where the input is a short narrative description and the output is a cinematic plan. We construct a prompt set covering diverse genres and scene types (e.g., dialogue-heavy, action-centric, interior/exterior, day/night). All methods are evaluated on the same prompt set to ensure comparability.

**Language Model Backbones.** All planning-based methods are instantiated with open-source instruction-tuned language models. We report results under three representative 7B-scale backbones: Qwen-7B-Instruct, LLaMA-2-7B-Chat, and Mistral-7B-Instruct. For each backbone, all compared methods and ablation variants use the same model to control for capacity.

### Compared Methods and Output Alignment.

We compare three representative paradigms: (i) prompt-based text-to-video planning, (ii) workflow-based planning with node-level structure, and (iii) free-form LLM planning without explicit schema constraints. Since our focus is planning quality rather than rendered videos, we evaluate the generated *plans*. For baselines that produce free-form prompts or workflow descriptions, we apply a deterministic field-mapping to extract comparable shot-level attributes when possible; otherwise, missing attributes are counted as missing fields by definition.

### 4.2 Automatic Evaluation: Structural Validity and Executability

We first evaluate whether generated outputs satisfy the requirements of visual executability. This evaluation is fully automatic and independent of human judgment.

**Metrics.** We report the following metrics:

- **Schema Validity Rate (%)**: The percentage of outputs that fully conform to the predefined schema, including required fields, data types, and admissible values.
- **Missing Field Rate (%)**: The proportion of shots that omit at least one required attribute.
- **Atomicity Violation Rate (%)**: The proportion of shots that contain multiple visual foci, actions, or camera operations within a single shot.

**Operational Definition and Examples.** An output is considered schema-valid if and only if it fully conforms to the predefined cinematic schema. For example, a shot specification that includes all required fields such as shot type, composition, camera movement, and focus, with correct types and values, is counted as valid. In contrast, free-form textual descriptions or partially structured outputs that omit required attributes or mix multiple actions in a single shot are counted as invalid.

Missing fields are detected by checking the absence of any required attribute in the schema. Atomicity violations occur when a single shot description combines multiple actions or camera movements that would correspond to multiple cinematic shots in professional practice (e.g., “the character enters the room, the camera pans, and another character reacts”).

441	<b>Results.</b> Table 1 summarizes the structural validity results across different planning paradigms.	483
442		484
443	We note that schema validity reaches 100%	485
444	for our method across all backbones because invalid outputs are prevented <i>by construction</i> via	486
445	schema-guided decoding rather than filtered post hoc. Therefore, differences across backbones	487
446	mainly manifest in higher-level content choices (e.g., phrasing or stylistic preferences), while structural	488
447	validity remains guaranteed.	489
448		490
449		491
450		492
451	<b>4.3 Human Evaluation: Cinematic Plan</b>	493
452	<b>Quality</b>	494
453	To assess the practical usefulness of generated cinematic plans, we conduct a human evaluation study	495
454	focusing on real-world creative workflows.	496
455		497
456	<b>Participants.</b> We recruit two groups of evaluators:	498
457		499
458	• Professionals from an AI short-drama production company.	500
459		501
460	• Graduate students majoring in film directing and cinematography.	502
461		503
462	All participants have prior experience with cinematic planning or storyboard creation.	504
463		505
464	<b>Protocol and Statistics.</b> We conduct a randomized, anonymized evaluation in which each participant rates outputs from all methods for the same set	506
465	of inputs. To reduce order effects, we randomize the presentation order per input. We report mean	507
466	Likert scores aggregated across participants and prompts.	508
467		509
468		510
469		511
470		512
471	<b>Evaluation Criteria.</b> Evaluators rate each method on a 5-point Likert scale (1: very poor, 5: excellent) along three dimensions:	513
472		514
473		515
474	• <b>Visual Clarity:</b> How clearly the plan specifies visual content and camera intent.	516
475		517
476	• <b>Shot Coherence:</b> How well shots are temporally and semantically organized.	518
477		519
478	• <b>Direct Usability:</b> How easily the plan can be directly used in storyboarding or video generation pipelines with minimal modification.	520
479		521
480		522
481	<b>Results.</b> Human evaluation results are reported in Table 2.	523
482		524
		525
		526
		527
		528
		529
		530
		531
		532
		533
		534
		535
		536
		537
		538
		539
		540
		541
		542
		543
		544
		545
		546
		547
		548
		549
		550
		551
		552
		553
		554
		555
		556
		557
		558
		559
		560
		561
		562
		563
		564
		565
		566
		567
		568
		569
		570
		571
		572
		573
		574
		575
		576
		577
		578
		579
		580
		581
		582
		583
		584
		585
		586
		587
		588
		589
		590
		591
		592
		593
		594
		595
		596
		597
		598
		599
		600
		601
		602
		603
		604
		605
		606
		607
		608
		609
		610
		611
		612
		613
		614
		615
		616
		617
		618
		619
		620
		621
		622
		623
		624
		625
		626
		627
		628
		629
		630
		631
		632
		633
		634
		635
		636
		637
		638
		639
		640
		641
		642
		643
		644
		645
		646
		647
		648
		649
		650
		651
		652
		653
		654
		655
		656
		657
		658
		659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700
		701
		702
		703
		704
		705
		706
		707
		708
		709
		710
		711
		712
		713
		714
		715
		716
		717
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800

Method	Qwen-7B-Instruct			LLaMA-2-7B-Chat			Mistral-7B-Instruct		
	SV	MF	AV	SV	MF	AV	SV	MF	AV
Prompt-based T2V	2.1	48.6	61.3	1.7	51.2	64.8	2.4	46.9	59.7
Workflow-based T2V	6.4	27.9	34.5	5.8	29.4	36.1	6.9	26.7	33.2
LLM Planning (Free-form)	4.8	36.2	42.7	4.1	38.9	45.3	5.2	34.6	40.8
<b>Ours (Visual Executability)</b>	<b>100.0</b>	<b>0.0</b>	<b>3.1</b>	<b>100.0</b>	<b>0.0</b>	<b>4.2</b>	<b>100.0</b>	<b>0.0</b>	<b>3.8</b>

Table 1: Automatic evaluation across different language model backbones. SV: Schema Validity (%), MF: Missing Fields (%), AV: Atomicity Violations (%).

Method	Qwen-7B-Instruct			LLaMA-2-7B-Chat			Mistral-7B-Instruct		
	VC	SC	DU	VC	SC	DU	VC	SC	DU
Prompt-based T2V	2.7	2.4	1.9	2.5	2.2	1.8	2.8	2.5	2.0
Workflow-based Planning	3.4	3.6	3.2	3.2	3.4	3.0	3.5	3.7	3.3
LLM Planning (No Schema)	3.1	2.9	2.5	2.9	2.7	2.3	3.2	3.0	2.6
<b>Ours (Visual Executability)</b>	<b>4.3</b>	<b>4.5</b>	<b>4.6</b>	<b>4.1</b>	<b>4.3</b>	<b>4.4</b>	<b>4.2</b>	<b>4.4</b>	<b>4.5</b>

Table 2: Human evaluation results across different language model backbones (5-point Likert scale). VC: Visual Clarity, SC: Shot Coherence, DU: Direct Usability.

**Discussion.** The ablation results in Table 3 reveal two complementary factors behind visual executability. First, removing shot-level decomposition (Scene-Level Only) substantially reduces shot consistency and usability, indicating that coarse scene descriptions do not provide sufficient granularity for execution: evaluators reported ambiguous shot boundaries and unclear camera intent, which prevents direct storyboard translation. Second, keeping shot-level structure but removing schema constraints (Shot-Level, No Schema) improves consistency compared to scene-level outputs, yet still yields structurally unstable plans with missing or inconsistent attributes, resulting in lower usability.

The full model combines shot-level granularity with hard schema constraints and achieves the best performance under all backbones. This confirms that visual executability is not obtained by adding format alone; rather, it emerges from jointly enforcing (i) an execution-oriented unit of planning (shots) and (ii) a verifiable interface (schema) that eliminates underspecification.

## 5 Qualitative Analysis

We present a qualitative example to illustrate how visual executability enables the generation of concrete and controllable cinematic plans, and how these plans translate into executable video outputs.

### 5.1 Case Study: From Story to Executable Video

**Input Story.** “A man waits alone at a bus stop late at night. Rain begins to fall. He checks his phone repeatedly, growing more anxious as time passes.”

**Generated Cinematic Plan.** Using our method, the language model produces the following shot-level plan (abridged for brevity):

Shot 1:	566
Shot Type: Wide shot	567
Composition: Rule of thirds	568
Camera Movement: Static	569
Focus: Empty bus stop under streetlight, light rain	570
	571
	572
Shot 2:	573
Shot Type: Medium shot	574
Composition: Centered	575
Camera Movement: Slow dolly in	576
Focus: Man standing alone, checking his phone	577
	578
Shot 3:	579
Shot Type: Close-up	580
Composition: Tight framing	581
Camera Movement: Static	582
Focus: Man's face, rain droplets visible, anxious expression	583
	584

The plan is fully schema-valid and specifies shot boundaries, camera intent, and visual focus for each shot, enabling direct execution without additional interpretation.

**Executable Image Output.** To qualitatively demonstrate visual executability, we instantiate the generated cinematic plan using representative image outputs rather than full video generation. Specifically, we render one image per selected shot using an open-source text-to-image diffusion model conditioned on the corresponding shot specification.

Figure 2 visualizes three shots generated from a single story input. Each image corresponds to one shot in the plan and reflects the intended shot type, composition, and visual focus. This example illustrates that the proposed shot-level schema produces

Variant	Qwen-7B-Instruct			LLaMA-2-7B-Chat			Mistral-7B-Instruct		
	SV	HU	SC	SV	HU	SC	SV	HU	SC
Scene-Level Only	1.3	2.4	1.9	0.9	2.2	1.7	1.5	2.5	2.0
Shot-Level, No Schema	5.1	3.2	3.4	4.6	3.0	3.2	5.4	3.3	3.5
<b>Shot-Level + Schema (Full)</b>	<b>100.0</b>	<b>4.6</b>	<b>4.5</b>	<b>100.0</b>	<b>4.4</b>	<b>4.3</b>	<b>100.0</b>	<b>4.5</b>	<b>4.4</b>

Table 3: Ablation study across different language model backbones. SV: Schema Validity (%), HU: Human Usability, SC: Shot Consistency.

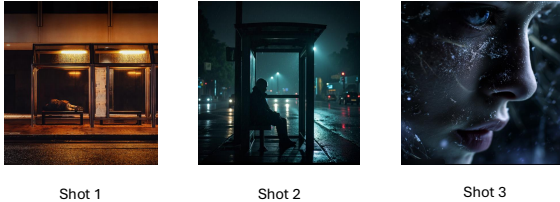


Figure 2: Executable image outputs corresponding to three shots in a generated cinematic plan. From left to right: (Shot 1) establishing wide shot of the environment, (Shot 2) medium shot focusing on the character in context, and (Shot 3) close-up emphasizing emotional detail. Each image is rendered independently from its shot specification, demonstrating the executability of the plan at the shot level.

based, and free-form planning baselines. Importantly, these improvements are consistent across multiple open-source language model backbones, indicating that executability arises from representation and constraint design rather than model-specific behavior.

We believe that visual executability provides a principled interface between language models and visual generation systems. Beyond cinematic planning, this paradigm opens up new opportunities for controllable multimodal generation, structured creative assistance, and human-AI collaboration in visual storytelling and content creation.

plans that can be directly translated into concrete visual realizations without additional interpretation or manual refinement.

## 5.2 Summary of Experimental Findings

Across automatic evaluation, human studies, ablation experiments, and backbone robustness analysis, results consistently demonstrate the effectiveness of enforcing visual executability. Shot-level schema constraints lead to structurally valid, coherent, and practically usable cinematic plans. These findings validate visual executability as a core design principle for language-driven cinematic planning systems.

## 6 Conclusion

We introduced visual executability as a design principle for language-driven cinematic planning. By representing plans at the shot level and enforcing schema-guided constraints during generation, our approach produces cinematic plans that are structurally valid, semantically coherent, and directly executable in real-world creative workflows.

Through automatic evaluation, human studies, ablation experiments, and qualitative analysis, we show that visual executability substantially improves both the reliability and usability of generated plans compared to prompt-based, workflow-

## 641 Limitations

642 While our results demonstrate the effectiveness of  
643 visual executability for cinematic planning, several  
644 limitations remain.

645 First, our framework relies on a predefined cin-  
646 ematic schema that encodes common shot-level  
647 attributes such as shot type, composition, and cam-  
648 era movement. Although this schema covers a wide  
649 range of standard cinematic practices, it may not  
650 fully capture highly unconventional or experimen-  
651 tal filmmaking styles. Extending the schema to  
652 support richer stylistic variation or alternative cine-  
653 matic grammars is an important direction for future  
654 work.

655 Second, while schema-guided decoding guaran-  
656 tees structural validity, it does not by itself ensure  
657 aesthetic quality or creative originality. As illus-  
658 trated by our qualitative examples, generated plans  
659 are executable and well-structured, but artistic ap-  
660 peal still depends on downstream rendering models  
661 and subjective human preferences. Balancing strict  
662 executability constraints with creative flexibility  
663 remains an open challenge.

664 Third, our human evaluation focuses on usabil-  
665 ity and clarity in professional workflows, rather  
666 than long-term narrative quality or audience recep-  
667 tion. Evaluating how executable cinematic plans  
668 influence end-user engagement or storytelling ef-  
669 fectiveness over longer video sequences is beyond  
670 the scope of this work.

671 Finally, although we evaluate our method across  
672 multiple open-source language model backbones,  
673 all experiments are conducted at a similar model  
674 scale. Investigating how visual executability inter-  
675 acts with larger models or multimodal foundation  
676 models is left for future research.

## 677 References

678 Hadj Batatia and Ilia Svetlichnyi. 2025. Feature execu-  
679 tion graphs: A human-ai co-programming paradigm  
680 for graph-driven llm code synthesis. In *Proceedings  
681 of the AAI Symposium Series*, volume 6, pages 184–  
682 191.

683 Liang Chen, Shuai Bai, Wenhao Chai, Weichu Xie,  
684 Haozhe Zhao, Leon Vinci, Junyang Lin, and Baobao  
685 Chang. 2025. Multimodal representation align-  
686 ment for image generation: Text-image interleaved  
687 control is easier than you think. *arXiv preprint  
688 arXiv:2502.20172*.

689 John Dagdelen, Alexander Dunn, Sanghoon Lee,  
690 Nicholas Walker, Andrew S Rosen, Gerbrand Ceder,

Kristin A Persson, and Anubhav Jain. 2024. Struc-  
691 tured information extraction from scientific text with  
692 large language models. *Nature communications*,  
693 15(1):1418. 694

Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani,  
695 Arjun Akula, Xuehai He, Sugato Basu, Xin Eric  
696 Wang, and William Yang Wang. 2023. Layoutgpt:  
697 Compositional visual planning and generation with  
698 large language models. *Advances in Neural Informa-  
699 tion Processing Systems*, 36:18225–18250. 700

Saibo Geng, Hudson Cooper, Michał Moskal, Samuel  
701 Jenkins, Julian Berman, Nathan Ranchin, Robert  
702 West, Eric Horvitz, and Harsha Nori. Json-  
703 schemabench: Evaluating constrained decoding with  
704 llms on efficiency, coverage and quality. 705

Maitrey Gramopadhye and Daniel Szafrir. 2023. Gener-  
706 ating executable action plans with environmentally-  
707 aware language models. In *2023 IEEE/RSJ Interna-  
708 tional Conference on Intelligent Robots and Systems  
709 (IROS)*, pages 3568–3575. IEEE. 710

Lin Guan, Karthik Valmeekam, Sarath Sreedharan,  
711 and Subbarao Kambhampati. 2023. Leveraging pre-  
712 trained large language models to construct and utilize  
713 world models for model-based task planning. *Ad-  
714 vances in Neural Information Processing Systems*,  
715 36:79081–79094. 716

Yingqing He, Menghan Xia, Haoxin Chen, Xiaodong  
717 Cun, Yuan Gong, Jinbo Xing, Yong Zhang, Xin-  
718 tao Wang, Chao Weng, Ying Shan, and 1 oth-  
719 ers. 2023. Animate-a-story: Storytelling with  
720 retrieval-augmented video generation. *arXiv preprint  
721 arXiv:2307.06940*. 722

Rong Hu, Ye Yang, Sen Liu, Zuchen Li, Jingyi Liu,  
723 Xingchen Ding, Hanchi Sun, and Lingli Ren. 2025.  
724 Large language model driven transferable key infor-  
725 mation extraction mechanism for nonstandardized  
726 tables. *Scientific Reports*, 15(1):29802. 727

Yu Meng, Jiaxin Huang, Yu Zhang, and Jiawei Han.  
728 2022. Generating training data with language mod-  
729 els: Towards zero-shot language understanding. *Ad-  
730 vances in Neural Information Processing Systems*,  
731 35:462–477. 732

James Oswald, Kavitha Srinivas, Harsha Kokel, Junkyu  
733 Lee, Michael Katz, and Shirin Sohrabi. 2024. Large  
734 language models as planning domain generators. In  
735 *Proceedings of the International Conference on Au-  
736 tomated Planning and Scheduling*, volume 34, pages  
737 423–431. 738

Vishal Pallagani, Bharath Chandra Muppasani, Kaushik  
739 Roy, Francesco Fabiano, Andrea Loreggia, Keerthi-  
740 ram Murugesan, Biplav Srivastava, Francesca Rossi,  
741 Lior Horesh, and Amit Sheth. 2024. On the prospects  
742 of incorporating large language models (llms) in au-  
743 tomated planning and scheduling (aps). In *Proceedings  
744 of the International Conference on Automated Plan-  
745 ning and Scheduling*, volume 34, pages 432–444. 746

747 Shreyas Sundara Raman, Vanya Cohen, Eric Rosen,  
748 Ifrah Idrees, David Paulius, and Stefanie Tellex. 2022.  
749 Planning with large language models via corrective  
750 re-prompting. In *NeurIPS 2022 Foundation Models  
751 for Decision Making Workshop*.

752 Karthik Valmeekam, Matthew Marquez, Sarath Sreed-  
753 haran, and Subbarao Kambhampati. 2023. On the  
754 planning abilities of large language models—a criti-  
755 cal investigation. *Advances in Neural Information  
756 Processing Systems*, 36:75993–76005.

757 Bailin Wang, Zi Wang, Xuezhi Wang, Yuan Cao, Rif  
758 A Saurous, and Yoon Kim. 2023a. Grammar prompt-  
759 ing for domain-specific language generation with  
760 large language models. *Advances in Neural Informa-  
761 tion Processing Systems*, 36:65030–65055.

762 Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Ji-  
763 ashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan  
764 Chen, Yaohui Wang, and 1 others. 2023b. In-  
765 ternvid: A large-scale video-text dataset for multi-  
766 modal understanding and generation. *arXiv preprint  
767 arXiv:2307.06942*.

768 Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu,  
769 Xiaojian Shawn Ma, and Yitao Liang. 2023c. De-  
770 scribe, explain, plan and select: interactive planning  
771 with llms enables open-world multi-task agents. *Ad-  
772 vances in Neural Information Processing Systems*,  
773 36:34153–34189.

774 Zhenyu Wu, Ziwei Wang, Xiuwei Xu, Jiwen Lu,  
775 and Haibin Yan. 2023. Embodied task plan-  
776 ning with large language models. *arXiv preprint  
777 arXiv:2307.01848*.

778 Dingyi Yang, Chunru Zhan, Ziheng Wang, Biao Wang,  
779 Tiezheng Ge, Bo Zheng, and Qin Jin. 2024. Syn-  
780 chronized video storytelling: Generating video nar-  
781 rations with structured storyline. *arXiv preprint  
782 arXiv:2405.14040*.

783 Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie,  
784 Chuyue Livia Sun, Jeff Huang, Cody Hao Yu, Shiyi  
785 Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonza-  
786 lez, and 1 others. 2024. Sglang: Efficient execution  
787 of structured language model programs. *Advances  
788 in neural information processing systems*, 37:62557–  
789 62583.

790 Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman,  
791 Haohan Wang, and Yu-Xiong Wang. 2023a. Lan-  
792 guage agent tree search unifies reasoning acting  
793 and planning in language models. *arXiv preprint  
794 arXiv:2310.04406*.

795 Wangchunshu Zhou, Yuchen Eleanor Jiang, Ethan  
796 Wilcox, Ryan Cotterell, and Mrinmaya Sachan.  
797 2023b. Controlled text generation with natural lan-  
798 guage instructions. In *International Conference on  
799 Machine Learning*, pages 42602–42613. PMLR.

800 Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi  
801 Feng, and Qibin Hou. 2024. Storydiffusion: Con-  
802 sistent self-attention for long-range image and video

generation. *Advances in Neural Information Process-  
ing Systems*, 37:110315–110340.

803  
804