

# CODE2VIDEO: A CODE-CENTRIC PARADIGM FOR EDUCATIONAL VIDEO GENERATION

Yanzhe Chen<sup>\*</sup> Kevin Qinghong Lin<sup>\*</sup> Mike Zheng Shou<sup>✉</sup>

Show Lab, National University of Singapore

<https://showlab.github.io/Code2Video/>

## ABSTRACT

While recent generative models advance pixel-space video synthesis, they remain limited in producing professional educational videos, which demand disciplinary knowledge, precise visual structures, and coherent transitions, limiting their applicability in educational scenarios. Intuitively, such requirements are better addressed through the manipulation of a renderable environment, which can be explicitly controlled via logical commands (*e.g.*, code). In this work, we propose **Code2Video**, a code-centric agent framework for generating educational videos via executable Python code. The framework comprises three collaborative agents: (*i*) *Planner*, which structures lecture content into temporally coherent flows and prepares corresponding visual assets; (*ii*) *Coder*, which converts structured instructions into executable Python codes while incorporating scope-guided auto-fix to enhance efficiency; and (*iii*) *Critic*, which leverages vision-language models (VLM) with visual anchor prompts to refine spatial layout and ensure clarity. To support systematic evaluation, we build **MMMC**, a benchmark of professionally produced, discipline-specific educational videos. We evaluate **MMMC** across diverse dimensions, including VLM-as-a-Judge aesthetic scores, code efficiency, and particularly, **TeachQuiz**, a novel end-to-end metric that quantifies how well a VLM, after unlearning, can recover knowledge by watching the generated videos. Our results demonstrate the potential of **Code2Video** as a scalable, interpretable, and controllable approach, achieving 40% improvement over direct code generation and producing videos comparable to human-crafted tutorials. The code and datasets are available at <https://github.com/showlab/Code2Video>.

## 1 INTRODUCTION

*“If you want to master something, teach it.” – Richard Feynman*

Recent advances in natural video generation have made remarkable progress in *pixel* space. End-to-end solutions, including diffusion-based (Ho et al., 2022a; Weng et al., 2024b) and autoregressive architectures (Weng et al., 2024a; Yuan et al., 2025), can synthesize visually compelling videos directly from text prompts (*i.e.*, **Text2Video**), achieving fine appearance and short-form fidelity. Yet these models struggle when the task requires long-form reasoning or multi-entity interaction (Li et al., 2024a). To overcome these limitations, recent works have moved toward multi-agent pipelines, where complex video generation is decomposed into collaborative subtasks, allowing iterative refinement, temporal structuring (Yuan et al., 2024; Huang et al., 2024; Xie et al., 2024).

Educational videos that aim to teach subject-specific knowledge face unique challenges in the reasoning era. Unlike short-form entertainment, educational content must integrate deep domain expertise (Clark & Mayer, 2023), carefully designed animations or transitions, and step-by-step reasoning (Bao et al., 2009) to support actual skill acquisition. This raises two fundamental challenges: (**i**) How to create high-quality educational videos that maintain both temporal coherence—concepts introduced, expanded, and reinforced in logical sequence—and spatial clarity—elements arranged legibly without occlusion; and (**ii**) How to evaluate educational videos beyond appearance, ensuring that they are educationally effective and semantically aligned with the intended learning topic. Existing video generation pipelines rarely satisfy these requirements, leaving a critical gap for agentic methods that unify temporal planning, spatial organization, and educational assessment.

<sup>\*</sup>Equal contribution

<sup>✉</sup>Corresponding author: mike.zheng.shou@gmail.com

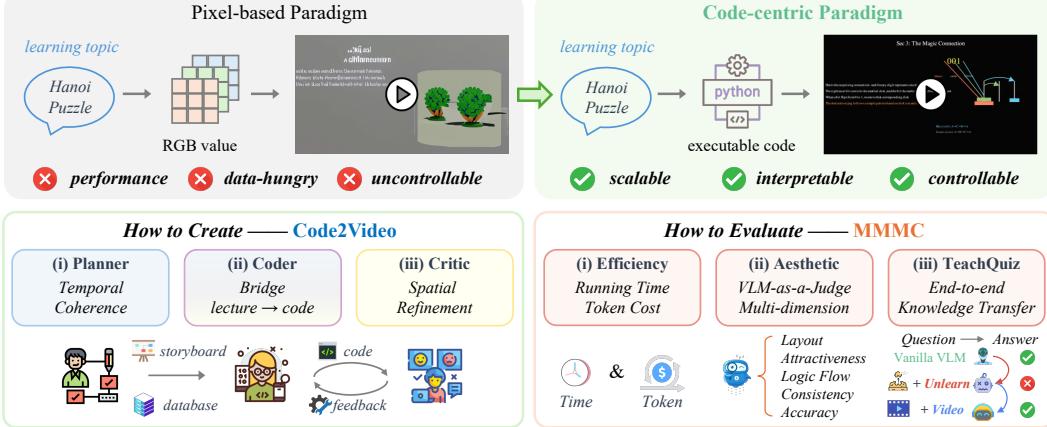


Figure 1: Overview of **Code2Video**. A code-centric paradigm for educational video generation, where Planner ensures temporal flow, Coder bridges instructions to executable animations, and Critic refines spatial layout. Evaluation is performed on **MMMC** with multi-dimensional metrics.

We are motivated by the intuition that code provides a uniquely suitable substrate for educational video generation. Unlike black-box models, code-centric pipelines are *scalable*, since new visualizations and external assets can be modularly integrated; *interpretable*, as every sequence, layout, and rendering decision is explicitly scripted and thus auditable; and *controllable*, enabling precise temporal sequencing and spatial organization through programmatic specification.

Building on these insights, we propose **Code2Video**, an agentic, code-centric framework for generating high-quality educational videos. The system decomposes the task into three agents: the *Planner* sequences concepts, examples, and recaps into a coherent lecture flow; the *Coder* translates structured instructions into executable Manim code, yielding precise, editable visualizations with consistent layout and timing; and the *Critic* leverages multimodal feedback and visual anchor prompts to refine spatial organization and ensure alignment with learning objectives. This tri-agent design explicitly models the temporal and spatial structure of instruction, while grounding the entire pipeline in transparent, reproducible, and extensible code.

To evaluate this paradigm, we propose **MMMC**, a benchmark reflecting the distinct goal of educational videos: teaching new knowledge. It comprises professionally produced, discipline-specific Manim tutorials across 13 domains (*e.g.*, topology, physics). Evaluation covers three complementary dimensions: (i) VLM-as-a-Judge aesthetic and structural quality; (ii) code efficiency, measuring generation time and token consumption; and (iii) **TeachQuiz**, a novel end-to-end knowledge-transfer metric that enforces unlearning of the target concept in a VLM, and then measures how effectively the generated video restores it. This multi-dimensional protocol directly probes educational efficacy and grounds a code-centric paradigm for video generation. Our results reveal clear trends: pixel-based models struggle with fine details and coherence, while direct code-centric generation improves TeachQuiz by 30%. Our full Planner–Coder–Critic pipeline further delivers a stable 40% gain. In human studies on TeachQuiz scores, agentically generated videos even outperform professional human-made tutorials, underscoring the *effectiveness of code-centric, agentic generation*.

Our contributions are summarized as follows:

- **A New Paradigm for Video Generation.** We introduce a new code-centric paradigm for educational video generation, positioning executable code as the unifying medium for temporal sequencing and spatial organization.
- **Effective Designs for Visual Animation Agent.** We highlight a modular agent design with three key components: (i) Planner expands an external database for reference, enabling parallel yet consistent storyboard; (ii) Coder ensures compilable code via automatic debugging and scope-guided repair; (iii) Critic refines spatial layout and clarity using visual anchor prompts.
- **New Benchmark with Well-designed Evaluation.** We present **MMMC**, the first benchmark for code-centric educational video generation with multi-dimensional evaluation of efficiency, aesthetics, and end-to-end knowledge transfer.

## 2 RELATED WORK

### 2.1 VIDEO GENERATION

Early text-to-video generation methods **(i)** extend diffusion models into the temporal domain via space-time UNets and latent 3D VAEs (Weng et al., 2024b; Ho et al., 2022b), achieving strong perceptual fidelity and longer durations (Yang et al., 2024; Li et al., 2024a; Xing et al., 2024). However, their reliance on *pixel-space* synthesis limits controllability, making precise layout and symbolic alignment—both critical for educational videos—difficult to realize. Autoregressive and progressive schemes (Li et al., 2024b; Gu et al., 2025; Wang et al., 2024; Xie et al., 2025) have improved long-form generation (Lu et al., 2024; Zhou et al., 2024), yet still struggle with board-like composition and stepwise exposition required in educational contexts (Li et al., 2024a; Liu et al., 2024). **(ii)** Recent advances in **multi-agent collaboration** show that decomposing tasks, coordinating tool use, and enabling iterative self-improvement can substantially enhance reasoning and generation (Yuan et al., 2024; Hu et al., 2024; Xie et al., 2024; Shen et al., 2024). While multi-agent frameworks have proven effective in domains such as web interaction, their application to video generation remains unexplored (Ku et al., 2025; Wu et al., 2024b). **(iii)** Building on this paradigm, we propose a **code-centric animation framework** for educational video synthesis. By elevating executable code as the generative substrate, our approach achieves symbolic layout, temporally structured exposition, and deterministic reproducibility—capabilities unattainable with pixel-level diffusion.

### 2.2 CODING AGENTS

Recent advances in LLM-based tool use have shown that agents can autonomously call APIs, retrieve specifications, and verify outputs, enabling neuro-symbolic modularity and robust task decomposition (Yao et al., 2023; Wang et al., 2025). By integrating code execution and tool invocation, representative methods extend language models beyond **text-only** reasoning, supporting complex workflows and project-level code generation (Patil et al., 2024; Liu et al., 2025; Gupta et al., 2024). Such developments demonstrate the potential of LLM agents to coordinate external retrieval, maintain memory across parallel processes, and incorporate feedback loops for iterative refinement (Li, 2025; Xu et al., 2025; Zhang et al., 2024). In parallel, research at the intersection of coding and visual reasoning shows that generating and executing programs can yield structured perception and controllable rendering (Pang et al., 2025; Zhu et al., 2025; Lin et al., 2025). **Visual programming** and visual-to-code approaches leverage program synthesis for compositional reasoning and spatial arrangement, with benchmarks translating images or text into executable code for charts, plots, and graphical interfaces (Wu et al., 2024a; Zhao et al., 2025; Wei et al., 2025; Yen et al., 2025). While these works bridge symbolic and visual domains, they largely focus on *static* figures or localized visual tasks (Xing et al., 2025; Wen et al., 2024; Ye et al., 2025; Jain et al., 2025). We advance this line by integrating code generation and visual synthesis for *dynamic* educational **video creation**.

## 3 MMMC BENCHMARK

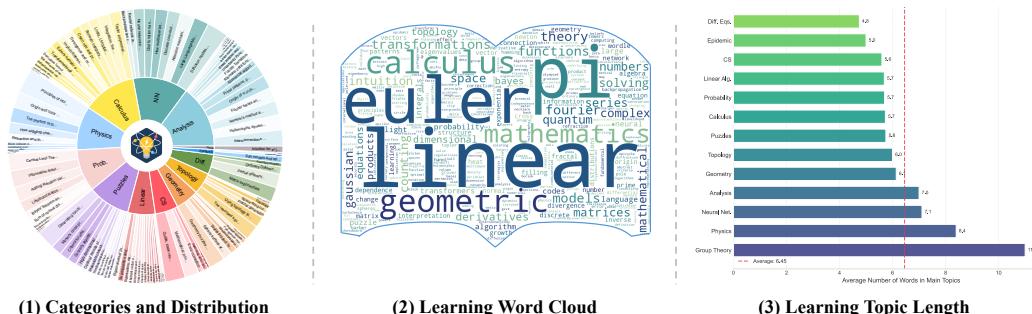


Figure 2: **MMMC overview**. (1) Left: distribution of 13 subject categories with exemplar learning topics; ring width encodes video duration. (2) Middle: learning topic word cloud highlighting core concepts. (3) Right: average learning topic length per category.

### 3.1 TASK FORMULATION

The task of code-centric educational video generation maps a learning query to executable *Manim* ([Manim Community Dev, 2025](#)) code whose rendering yields a tutorial video. The challenge lies in multi-step reasoning, precise temporal sequencing, and spatial coherence, where even minor syntax errors can nullify execution. We adopt *Manim* for its fine-grained spatiotemporal control, symbolic expressivity, and demonstrated effectiveness in expert-produced instructional videos.

### 3.2 DATA CURATION AND STATISTICS

We construct MMMC, a benchmark for code-driven educational video generation, under two criteria: (i) *educational relevance*—each learning topic is an established concept worth teaching; and (ii) *executable grounding*—each concept aligns with a high-quality Manim reference, ensuring practical realizability. We source from the complete 3Blue1Brown (3B1B) YouTube corpus, known for its instructional impact and expert Manim craftsmanship. After filtering out non-instructional items (e.g., Q&A), we curate 117 long-form videos spanning 13 subject areas, including *calculus*, *geometry*, *probability*, and *neural networks*. To enrich supervision, we segment videos using author-provided timestamps into 339 semantically coherent sub-clips, yielding 456 units in total. An LLM then extracts concise learning topics (avg. 6.3 words) from titles, descriptions, and metadata, producing a clean mapping from videos to educationally grounded units (details in §A.1.5). On average, a full-length video lasts 1014 seconds (~16.9 minutes), while a segmented clip spans 201 seconds (~3.35 minutes), thus balancing long-horizon reasoning with fine-grained supervision. Figure 2 visualizes topical diversity with a hierarchical donut plot: the inner ring denotes 13 categories, while the outer ring shows individual topics with arc width proportional to cumulative duration. This structure highlights both the breadth of coverage and the temporal richness of MMMC, establishing it as a challenging and representative benchmark for educational video generation.

### 3.3 EVALUATION METRICS

Unlike conventional video generation, educational videos are valued less for visual fidelity than for how effectively they convey knowledge. This makes standard synthesis metrics inadequate. We therefore design a three-pronged evaluation across **aesthetics**, **knowledge convey**, and **efficiency**:

**VLM-as-Judges.** Since human judgments of video quality are inherently subjective, we adopt a VLM-as-judges protocol ( $P_{\text{aesth}}$ ) to approximate user perception across five axes: (i) *Element Layout (EL)* — clarity and spatial arrangement of visual components. (ii) *Attractiveness (AT)* — overall engagement and ability to capture learners’ attention. (iii) *Logic Flow (LF)* — coherence in temporal presentation of concepts. (iv) *Visual Consistency (VC)* — stylistic stability across frames and sections. (v) *Accuracy & Depth (AD)* — correctness and richness of the presented knowledge. Each dimension is rated on a 100-point scale.

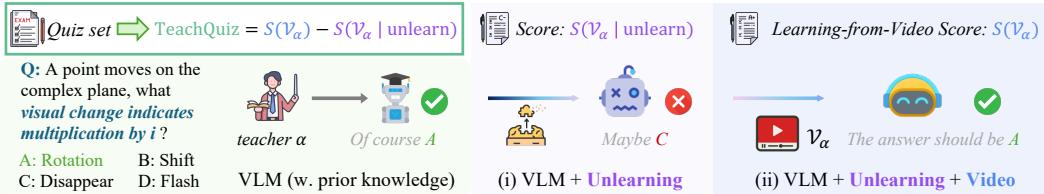


Figure 3: TeachQuiz: score gap between *Learning-from-Video* and *Unlearning* stages.

**TeachQuiz.** The goal of educational video generation is not merely visual plausibility, but effective knowledge transfer. To evaluate this, we introduce TeachQuiz, a two-stage protocol grounded in a quiz set  $\mathcal{Q}(\mathcal{K}) = (q_i, y_i)_{i=1}^N$  for a given concept  $\mathcal{K}$ , and  $Y$  denotes ground-truth answers. We consider multiple teachers  $\alpha, \beta$ , each producing a video  $\mathcal{V}_\alpha, \mathcal{V}_\beta$ . A student model  $\phi$  is tasked with watching the video and answering questions:

$$S(\mathcal{V}_\alpha) = \mathbf{1}[\phi(Q, \mathcal{V}_\alpha) = Y] \quad (1)$$

If  $S(\mathcal{V}_\alpha) > S(\mathcal{V}_\beta)$ , then teacher  $\alpha$  is the stronger instructor.

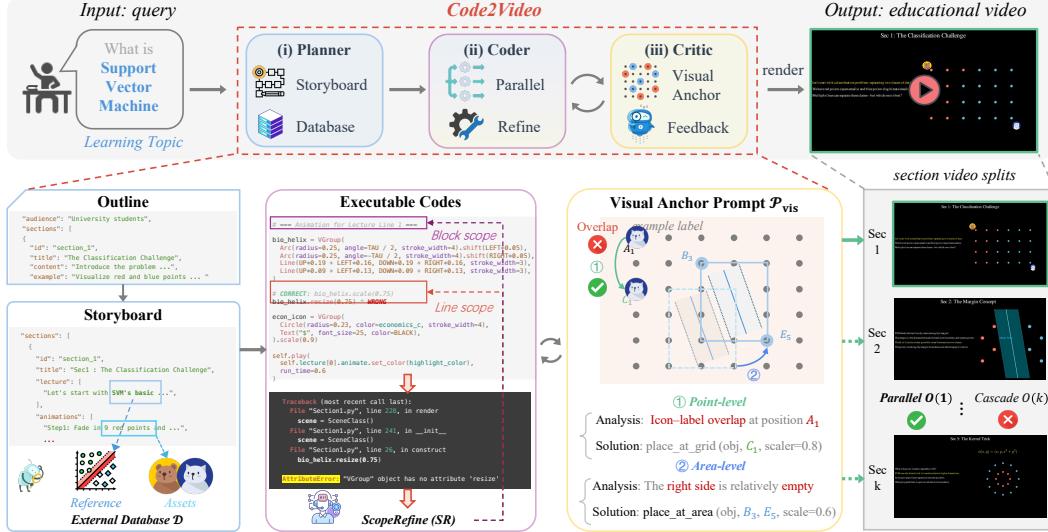


Figure 4: **Illustration of Code2Video.** Given a user inquiry, Code2Video aims to render an educational video via Manim code writing: (i) **the Planner** converts a learning topic into a storyboard and retrieves visual assets; (ii) **the Coder** performs parallel code synthesis with scope-guided refinement to ensure efficiency and temporal consistency; (iii) **the Critic** uses visual anchor prompts to iteratively adjust spatial layout and clarity, yielding reproducible, educationally structured videos.

However, a key challenge is that *many quiz items are already be learn by top-performing VLMs (i.e., answer correctly without watching the video)*. Thus, absolute accuracy alone does not measure teaching quality. Instead, a good educational video should improve knowledge acquisition relative to a controlled baseline. We enforce this through two steps: (i) **Unlearning**. Apply  $\mathcal{P}_{\text{unlearn}}$  to block prior access to  $\mathcal{K}$ , yielding a knowledge-removed baseline. (ii) **Learning-from-Video**. Expose the model to  $\mathcal{V}$  under  $\mathcal{P}_{\text{learn}}$ , testing whether the video itself enables recovery of the knowledge. The final *TeachQuiz* score measures relative improvement:

$$\tilde{S}(\mathcal{V}_\alpha) = S(\mathcal{V}_\alpha) - S(\mathcal{V}_\alpha | \text{unlearn}) \quad (2)$$

which isolates the contribution of the video by subtracting the unlearned baseline. Higher  $\tilde{S}$  indicates stronger knowledge transfer induced by the generated video.

**Token Cost and Generation Time.** Beyond output quality, an equally important dimension is how economically a model can generate effective videos. We measure efficiency by *average code generation time* and *token usage per video*, reflecting scalability and feasibility in large-scale or interactive educational settings where latency and resource costs are critical.

## 4 METHOD: CODE2VIDEO

**Overview.** As illustrated in Fig. 4, given a topic query  $\mathcal{Q}$ , Code2Video output a video  $\mathcal{V}$ , which consists of three stages: (i) **Planner** structures topics into storyboards with reference assets, (ii) **Coder** translates each section into executable Manim code using parallel synthesis and an effective debugging, and (iii) **Critic** refines rendered videos through a novel visual prompt and VideoLLM feedback to ensure spatial coherence and educational clarity.

### 4.1 PLANNER: QUERY TO STORYBOARD

Generating coherent educational videos requires careful organization of temporal structure. We design the Planner to decompose a topic  $\mathcal{Q}$  into two stages: outline generation for high-level ordering, and storyboard construction for stepwise realization. This preserves logical flow while capturing cross-section dependencies.

**(i) Outline Generation.** Given a topic  $\mathcal{Q}$ , the Planner produces an outline  $\mathcal{O} = o_1, \dots, o_n$ , where each  $o_i$  contains a unique identifier, section title, content summary, and illustrative examples. Crucially, the Planner also considers the intended audience (e.g., trigonometric functions for middle school, Fourier’s law for undergraduates), ensuring level-appropriate structure. Formally,  $\mathcal{O} \leftarrow \mathcal{P}_{\text{outline}}(\mathcal{Q})$ , where  $\mathcal{O} = \{o_1, \dots, o_n\}$  and each  $o_i$  encodes the section-level metadata and educational intent. By explicitly specifying audience and structure, the outline establishes the temporal skeleton for the subsequent video, guiding both pacing and sequencing.

**(ii) Storyboard Construction.** The second stage converts the outline  $o$  into a detailed storyboard  $s$ . Each section in  $s$  includes title, lecture lines, and corresponding animations, with  $s_i \leftarrow \mathcal{P}_{\text{storyboard}}(o_i)$ . The storyboard specifies the temporal sequence of lecture lines and paired animations, bridging high-level planning with concrete visual content.

**External Database.** To enhance factual accuracy and visual fidelity, the Planner integrates an external database  $\mathcal{D}$ . It includes (a) *reference images* aligned with the topic to anchor complex concepts and reduce hallucination, and (b) *visual assets* (e.g., icons, logos) that are difficult to generate from scratch. These assets  $\mathcal{A}$  are automatically identified via a prompt  $\mathcal{P}_{\text{asset}}$  analyzing the storyboard,  $a_i \leftarrow \mathcal{P}_{\text{asset}}(s_i)$ , and stored in a persistent cache  $\mathcal{D}_{\text{asset}}$ . Caching enables reuse across sections, preventing redundant generation and ensuring visual consistency. Please refer to § A.1.6 for more details and examples about  $\mathcal{D}$ .

## 4.2 CODER: STORYBOARDS TO EXECUTABLE CODE

The Coder  $\mathcal{G}$  translates each section of the storyboard  $s$  and the cached assets  $a$  into executable Manim code  $C = \{c_1, \dots, c_n\}$ , where each  $c_i$  corresponds to a storyboard  $s_i$ .

**(i) Parallel Code Generation.** A central bottleneck in full-code synthesis is generation time, as end-to-end Manim code production for a single educational video—including generation, debugging, and rendering—can exceed two hours. To address this bottleneck, we parallelize the pipeline by decoupling serial steps—code generation, debugging, and refinement—so that each section is synthesized and fixed independently. Each section is conditioned on its storyboard and shared assets  $\mathcal{A}$ :  $c_i = \mathcal{P}_{\text{coder}}(s_i, \mathcal{A})$ . Notably, asset sharing across sections ensures temporal consistency while retaining the efficiency benefits of parallelization.

**(ii) Effective Debugging.** Even strong LLMs rarely produce fully executable code in one pass. Naïve strategies that concatenate all code with the full error log are costly in both time and tokens. We propose **ScopeRefine (SR)**, a hierarchical, scope-guided repair strategy, as illustrated in Fig. 4 middle bottom: (a) *Line scope*: isolate the error line plus immediate context,  $S_1 = \text{line} \pm 1$ , attempt up to  $K_1$  local fixes. (b) *Block scope*: if unresolved, expand to the lecture-line block  $S_2 = \mathcal{B}_{i,j}$  with up to  $K_2$  repair attempts. (c) *Global scope*: as a last resort, regenerate the full section  $c_i$  from  $s_i$ . This progressive “*Go-to style*” repair minimizes token usage and latency while ensuring high reliability, effectively bridging parallel generation with robust debugging.

## 4.3 CRITIC: EFFECTIVE VISUAL REFINEMENT

Even after debugging ensures executability, the generated code may still yield unsatisfactory visual outcomes. LLMs and VLMs often fail to provide actionable feedback due to **limited spatial awareness** (Cheng et al., 2024; Zha et al., 2025). In practice, models can identify issues (e.g., “the cat icon is misplaced”) but struggle to provide actionable corrections. They often fail to indicate the direction or distance needed to adjust the element, which makes text-only refinement inadequate.

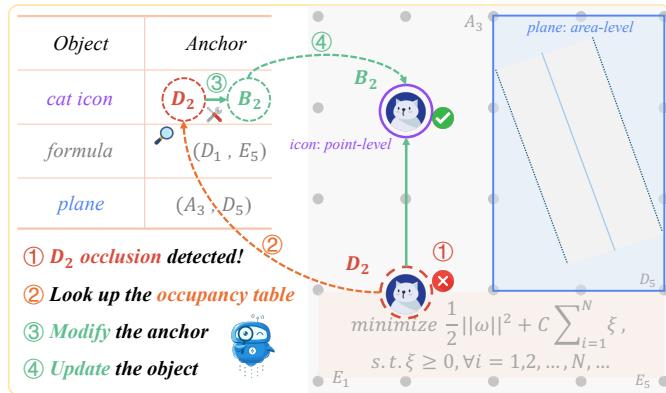


Figure 5: Illustration of visual anchor prompt ( $\mathcal{P}_{\text{vis}}$ ).

**(i) Visual Anchor Prompt ( $P_{\text{vis}}$ ).** We introduce  $P_{\text{vis}}$ , a textual prompt that discretizes the 2D canvas into a  $6 \times 6$  grid of predefined anchor points. Each grid cell is mapped to fixed Manim coordinates, allowing LLM-specified locations to be directly converted into executable code. Placement follows two granularities, as illustrated in Figure 5: *(a) point-level*, where small elements (e.g., symbols, short labels) occupy a single anchor; and *(b) region-level*, where larger elements are assigned to a bounding box spanning multiple anchors. This discretization transforms the placement task from a *continuous positioning problem* into a *discrete anchoring problem*, serving as a visual debugging “go-to”, which substantially reduces the difficulty for LLMs to produce valid layouts.

**(ii) VideoLLM for Code Feedback.** To detect violations and refine placement, the Critic inspects the rendered video  $\mathcal{V}_i$  alongside its section code  $c_i$ . During parallel code generation, we maintain an *occupancy table* that records each element’s assigned anchors (point or region), scaling factor, and corresponding code lines. This design serves two purposes: (a) it makes all assets indexable, allowing the Critic to quickly trace a visual issue back to its source code; and (b) it reveals available anchors, enabling conflict-free reallocation. With this structured view, the Critic efficiently detects three common issues: overlapping elements within a cell, lecture lines occluded by animations, and large unused regions creating visual imbalance. These findings are incorporated into a refinement prompt  $P_{\text{refine}}$ , yielding optimized code:  $\tilde{c}_i = P_{\text{refine}}(c_i, \mathcal{V}_i)$  and final video  $\tilde{\mathcal{V}} = \text{Render}(\{\tilde{c}_i\}_{i=1}^n)$ . By combining anchor-based guidance, indexable, occupancy-aware adjustment, and multimodal feedback, the Critic overcomes the limitations of text-only debugging.

## 5 EXPERIMENT

### 5.1 IMPLEMENTATION DETAILS

**Baselines.** We compare four types of approaches: ◇ *Human-crafted*, expert-designed Manim videos as an upper bound; ◇ *Pixel-based Diffusion*, text-to-video models: *OpenSora-v2* (Peng et al., 2025), *Wan2.2-T2V-A14B* (Wan et al., 2025), and *Veo3* (Google DeepMind, 2025); ◇ *CodeLLM Generation*, where an LLM directly generates Manim code from a learning topic; ◇ *Agentic Generation (ours)*, a Planner–Coder–Critic pipeline. We evaluate across diverse models: *Claude Opus 4.1* (Anthropic, 2025), *GPT-4o*, *GPT-o4 mini*, *GPT-4.1*, *GPT-5* (OpenAI, 2025), *Gemini-2.5 Pro* (Imran & Almusharraf, 2024), with *Gemini-2.5 Pro* serving as Critic for refinement. **Evaluation.** Aesthetics are judged by *Gemini-2.5 Pro* (VLM-as-a-Judge), and quantify knowledge transfer with TeachQuiz. **Resources.** Reference images are retrieved from Google Images, and visual assets from Iconfinder<sup>1</sup>. All prompts are documented in § A.2.

### 5.2 MAIN RESULTS

Table 1 compares Code2Video with human-crafted videos, pixel-based models, and code LLM baselines, evaluated on Efficiency, Aesthetics (AES), and knowledge transfer (TeachQuiz). Our analysis yields several insights: **(i) Pixel-based models underperform.** They obtain the lowest scores on both AES and TeachQuiz, particularly struggling with LF due to weak control over text grounding, animation timing, and cross-frame coherence. **(ii) Direct code-centric generation delivers clear improvements.** Rendering videos from LLM-produced Manim code outperforms pixel-based models, underscoring code as an effective medium for controllable and coherent educational video generation. **(iii) Our agentic framework delivers stable and consistent improvements.** Across different backbone LLMs, Code2Video achieves significant performance boosts. For instance, with Claude Opus 4.1, AES improves by 50% and TeachQuiz by 46%. These gains arise from distinct components: visual anchor points drive improvements in element layout, while the Planner enhances LF and AD. However, limitations remain in AT and VC, pointing to opportunities for refinement. **(iv) Human-made videos remain strong.** Although Code2Video narrows the gap, professional videos still lead in storytelling, nuanced sequencing, and explanatory depth. This highlights the next frontier: advancing agentic pipelines toward *professional-quality long educational videos*.

**Qualitative Analyses.** Figure 6 illustrates that our code-driven pipeline produces videos with clear text and formulas, stable layouts without occlusions, and stepwise alignment with lecture lines. In contrast, the pixel-based model (Veo3) often generates blurry or corrupted text, inconsistent styles,

---

<sup>1</sup><https://www.iconfinder.com>

Table 1: Results across Efficiency, Aesthetics, and TeachQuiz (Quiz). Efficiency: Time (**avg minutes** per topic) and Token (avg **token consumption** per topic). Aesthetics: Element Layout (EL), Attractiveness (AT), Logic Flow (LF), Visual Consistency (VC), Accuracy & Depth (AD).

Method	Efficiency (↓)		Aesthetics (↑)						Quiz (↑)
	Time	Token (K)	EL	AT	LF	VC	AD	Avg	
Human-made 3B1B	–	–	98.3	100	100	100	100	99.7	97.1
<i>Pixel-based Diffusion</i>									
OpenSora-v2	27.6	–	0.0	5.0	0.0	0.0	13.3	3.7	0.0
Wan2.2-T2V-A14B	17.4	–	0.0	10.0	0.0	0.0	20.0	6.0	0.0
Veo3	2.3	–	0.0	15.0	0.0	5.0	25.0	9.0	2.5
<i>Code LLM</i>									
GPT-5	1.8	1.1	27.0	28.0	28.0	54.5	26.0	32.7	36.5
GPT-4.1	2.1	1.2	30.5	34.5	39.0	42.0	24.8	34.2	37.0
Claude Opus 4.1	2.8	2.3	47.5	40.0	26.5	56.6	18.4	37.8	40.0
<i>Code2Video Agent (Ours)</i>									
Code2Video Gemini-2.5 Pro	15.5	41.8	70.3	60.3	44.3	37.6	74.7	57.4	72.0
Code2Video GPT-4o	14.1	32.7	70.3	58.3	54.6	48.5	68.3	60.0	44.0
Code2Video GPT-o4 mini	16.8	49.2	77.0	52.8	73.0	57.2	79.0	67.8	48.5
Code2Video GPT-5	8.8	19.3	75.5	60.5	81.8	63.6	79.7	72.2 <small>+39.5</small>	80.0 <small>+43.5</small>
Code2Video GPT-4.1	15.4	30.8	82.8	65.6	95.0	68.0	83.7	79.0 <small>+44.8</small>	82.0 <small>+45.0</small>
Code2Video Claude Opus 4.1	13.8	43.1	90.6	79.7	93.3	84.2	91.9	<b>87.9</b> <small>+50.1</small>	<b>86.0</b> <small>+46.0</small>

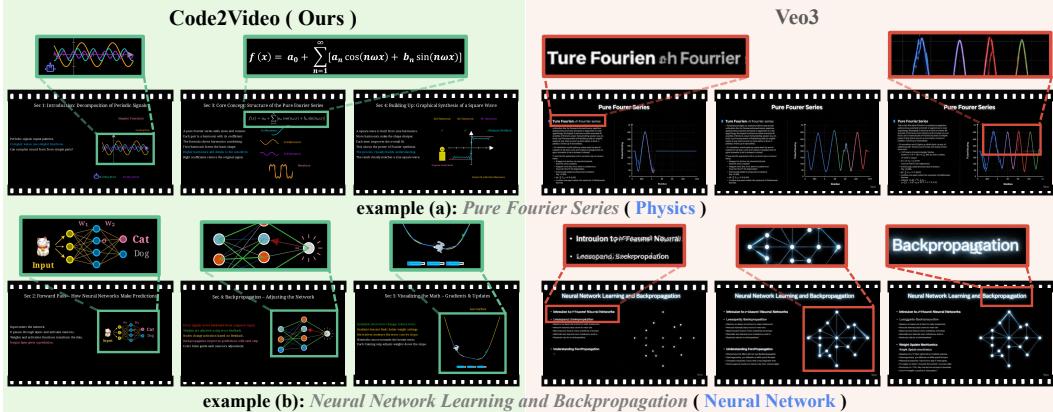


Figure 6: Qualitative comparison between *Code2Video* and *Veo3*. Our approach generates videos with coherent logic flow, consistent semantics, and interpretable layouts.

and drifting visuals, weakening semantic grounding. Overall, code-driven synthesis ensures better spatial stability and clearer knowledge presentation. Additional cases are provided in § A.1.7.

### 5.3 ABLATION STUDIES

**Effects by Individual Components.** Table 2 highlights several observations. First, TeachQuiz is more sensitive than Aesthetics, revealing *knowledge-transfer gaps even when videos still look visually acceptable*. Second, the Planner is essential: removing it collapses both metrics ( $\approx 41$  points), underscoring that high-level lecture planning and temporal sequencing are the backbone of effective teaching videos. Third, other modules provide complementary gains: the External Database improves conceptual grounding, Visual Anchors stabilize layouts, and the Critic ensures refinement—each modest alone, but jointly essential for robustness. These results **highlight that structured visual guidance and iterative refinement are crucial** for producing visually clear videos that effectively convey knowledge.

**Efficiency Components.** Table 3 evaluates efficiency-oriented modules. Removing parallel execution greatly increases latency ( $15.4 \rightarrow 86.6$  minutes). Without ScopeRefine (SR), we test two

Table 2: Effect of different components on quality: TeachQuiz / Aesthetics avg. score.

Method	Aesthetics	Quiz
Code2Video <code>chat-4.1</code> ( $\diamond$ )	<b>79.0</b>	<b>82.0</b>
◦ w/o Planner	38.1 $-40.9$	40.5 $-41.5$
◦ w/o External Database	68.1 $-10.9$	52.0 $-30.0$
◦ w/o Visual Anchor	69.2 $-9.8$	55.2 $-26.8$
◦ w/o Critic	72.5 $-6.5$	60.7 $-21.3$

Table 3: Effect of efficiency components: run-time avg. time / token consumption.

Method	Time (m)	Token (K)
Code2Video <code>chat-4.1</code> ( $\diamond$ )	<b>15.4</b>	<b>30.8</b>
◦ w/o parallel	86.6 $5.6\times$	30.8
◦ w/o SR $\rightarrow$ w. Retry	42.9 $2.8\times$	49.8 $1.6\times$
◦ w/o SR $\rightarrow$ w. Debug	39.2 $2.5\times$	42.1 $1.4\times$
◦ w/o parallel & SR	149.8 $9.7\times$	52.6 $1.7\times$

alternative debugging methods: (i) *Retry*, which regenerates the section upon any error; (ii) *Full-code Debug*, which feeds the entire code and error log to the LLM to regenerate the section. In both cases, error correction is costly, highlighting the importance of SR’s localized, scope-aware repair. Removing both mechanisms produces prohibitive overheads. These results underscore that parallel synthesis and scope-aware repair are essential for scalable, code-centric video generation.

Table 4: **Human study** on Aesthetics, TeachQuiz (Quiz), Completion Willingness (CW), and Average Ranking (AR). Results align with VLM-based trends but show sharper score contrast, lower tolerance for layout errors, and reduced engagement in longer-duration videos.

Method	Duration	Aesthetics ( $\uparrow$ )						Quiz ( $\uparrow$ )	CW ( $\uparrow$ )	AR ( $\downarrow$ )
		EL	AT	LF	VC	AD	Avg			
Human-made 3B1B	16.9 min	98.9	97.2	91.3	98.0	97.0	96.5	78.8	36.2	1.2
Pixel-based <code>veo3</code>	8.0 s	12.6	4.4	1.1	24.4	1.1	8.5	8.0	46.8	5.0
Code LLM <code>Claude Opus 4.1</code>	0.9 min	16.1	41.1	55.6	71.1	72.2	51.2	56.6	15.0	3.9
Code2Video <code>Gemini-2.5 Pro</code>	1.6 min	26.7	68.3	78.1	90.2	81.0	68.9	65.3	47.4	3.1
Code2Video <code>Claude Opus 4.1</code>	2.0 min	60.2	89.3	84.6	92.0	83.1	81.8	<b>80.3</b>	<b>64.0</b>	1.8

**Human Study Evaluation.** We conduct a five-group user study (6 middle school, 2 undergraduate volunteers per group), where each participant watches one video type and answers 5 quiz questions for 20 learning topics. We measure Completion Willingness (**CW**, proportion finishing the video before answering, max score is 100) and Average Ranking (**AR**, mean preference across video types, 1 is the best). Table 4 reveals four patterns: (i) **Clearer separation**. Human ratings follow the same overall trends as VLM-based scores but with stronger contrast: high-quality videos are rated in the upper range ( $> 90$ ), while low-quality videos cluster near the lower bound ( $< 10$ ). (ii) **Sensitivity to layout errors**. Participants give lower layout scores (EL) to videos from Code2Video, as humans are highly sensitive to even brief occlusions, whereas VideoLLMs often miss such frame-level issues. (iii) **Attention span limits**. Human attention is inherently limited: to perform well on the quiz, participants must follow the full flow of knowledge details in the video. This requires not only *strong logical coherence* and *engaging presentation* but also a *reasonable duration* that allows sustained high attention for effective knowledge absorption. (iv) **Strong consistency**. Aesthetics and TeachQuiz scores are strongly correlated( $r = 0.971, p = 0.0059$ ): visually appealing videos keep students engaged, leading to higher learning outcomes. Overall, the human study underscores that both structural clarity and visual appeal are decisive levers for learning efficacy, complementing the automated metrics. *Future work requires agent designs that explicitly account for human attention and patience, ensuring videos maintain fine-grained details while minimizing perceptual fatigue.*

## 6 CONCLUSION

We have introduced a novel, code-centric paradigm for educational video generation, establishing executable code as the unifying medium for both temporal sequencing and spatial organization. Building on this paradigm, our tri-agent architecture *Code2Video* enables controllable and interpretable generation with multimodal feedback. To systematically evaluate this paradigm, we introduce *MMMC*, targeting efficiency, aesthetics, and knowledge transfer. Together, our paradigm, architecture, and benchmark chart a clear path for future research on leveraging code as a medium for high-quality, structured, and interpretable educational content generation. Future work includes broadening the video scope and developing more lightweight, scalable agent frameworks.

## REFERENCES

- Anthropic. Claude opus 4.1. <https://www.anthropic.com/clause>, 2025.
- Lei Bao, Tianfan Cai, Kathy Koenig, Kai Fang, Jing Han, Jing Wang, Qing Liu, Lin Ding, Lili Cui, Ying Luo, et al. Learning and scientific reasoning. *Science*, 323(5914):586–587, 2009.
- An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093, 2024.
- Ruth C Clark and Richard E Mayer. *E-learning and the science of instruction: Proven guidelines for consumers and designers of multimedia learning*. John Wiley & Sons, 2023.
- Google DeepMind. Veo 3: Generative video model. <https://deepmind.google/technologies/veo/>, 2025.
- Yuchao Gu, Weijia Mao, and Mike Zheng Shou. Long-context autoregressive video modeling with next-frame prediction. *arXiv preprint arXiv:2503.19325*, 2025.
- Tanmay Gupta, Luca Weihs, and Aniruddha Kembhavi. Codenav: Beyond tool-use to using real-world codebases with llm agents. *arXiv preprint arXiv:2406.12276*, 2024.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022a.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in neural information processing systems*, 35:8633–8646, 2022b.
- Panwen Hu, Jin Jiang, Jianqi Chen, Mingfei Han, Shengcai Liao, Xiaojun Chang, and Xiaodan Liang. Storyagent: Customized storytelling video generation via multi-agent collaboration. *arXiv preprint arXiv:2411.04925*, 2024.
- Kaiyi Huang, Yukun Huang, Xuefei Ning, Zinan Lin, Yu Wang, and Xihui Liu. Genmac: compositional text-to-video generation with multi-agent collaboration. *arXiv preprint arXiv:2412.04440*, 2024.
- Muhammad Imran and Norah Almusharraf. Google gemini as a next generation ai educational tool: a review of emerging educational technology. *Smart Learning Environments*, 11(1):22, 2024.
- Vyoman Jain, Shiva Golugula, Motamarri Sai Sathvik, et al. Manimator: Transforming research papers into visual explanations. *arXiv preprint arXiv:2507.14306*, 2025.
- Max Ku, Thomas Chong, Jonathan Leung, Krish Shah, Alvin Yu, and Wenhua Chen. Theoremexplainagent: Towards video-based multimodal explanations for llm theorem understanding. *arXiv preprint arXiv:2502.19400*, 2025.
- Chengxuan Li, Di Huang, Zeyu Lu, Yang Xiao, Qingqi Pei, and Lei Bai. A survey on long video generation: Challenges, methods, and prospects. *arXiv preprint arXiv:2403.16407*, 2024a.
- Xinzhe Li. A review of prominent paradigms for llm-based agents: Tool use, planning (including rag), and feedback learning. In *Proceedings of the 31st International Conference on Computational Linguistics*, pp. 9760–9779, 2025.
- Zongyi Li, Shujie Hu, Shujie Liu, Long Zhou, Jeongsu Choi, Lingwei Meng, Xun Guo, Jinyu Li, Hefei Ling, and Furu Wei. Arlon: Boosting diffusion transformers with autoregressive models for long video generation. *arXiv preprint arXiv:2410.20502*, 2024b.
- Kevin Qinghong Lin, Linjie Li, Difei Gao, Zhengyuan Yang, Shiwei Wu, Zechen Bai, Stan Weixian Lei, Lijuan Wang, and Mike Zheng Shou. Showui: One vision-language-action model for gui visual agent. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 19498–19508, 2025.

Kaiyuan Liu, Youcheng Pan, Yang Xiang, Daojing He, Jing Li, Yexing Du, and Tianrun Gao. Projecteval: A benchmark for programming agents automated evaluation on project-level code generation. *arXiv preprint arXiv:2503.07010*, 2025.

Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024.

Yu Lu, Yuanzhi Liang, Linchao Zhu, and Yi Yang. Freelong: Training-free long video generation with spectralblend temporal attention. *Advances in Neural Information Processing Systems*, 37: 131434–131455, 2024.

Manim Community Dev. Manim community v0.19.0. <https://github.com/ManimCommunity/manim>, 2025.

OpenAI. Chatgpt-series. <https://openai.com>, 2025.

Wei Pang, Kevin Qinghong Lin, Xiangru Jian, Xi He, and Philip Torr. Paper2poster: Towards multimodal poster automation from scientific papers. *arXiv preprint arXiv:2505.21497*, 2025.

Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive apis. *Advances in Neural Information Processing Systems*, 37:126544–126565, 2024.

Xiangyu Peng, Zangwei Zheng, Chenhui Shen, Tom Young, Xinying Guo, Binluo Wang, Hang Xu, Hongxin Liu, Mingyan Jiang, Wenjun Li, et al. Open-sora 2.0: Training a commercial-level video generation model in 200 k. *arXiv preprint arXiv:2503.09642*, 2025.

Leixian Shen, Haotian Li, Yun Wang, and Huamin Qu. From data to story: Towards automatic animated data video creation with llm-based multi-agent systems. In *2024 IEEE VIS Workshop on Data Storytelling in an Era of Generative AI (GEN4DS)*, pp. 20–27. IEEE, 2024.

Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.

Hongru Wang, Cheng Qian, Manling Li, Jiahao Qiu, Boyang Xue, Mengdi Wang, Heng Ji, and Kam-Fai Wong. Toward a theory of agents as tool-use decision-makers. *arXiv preprint arXiv:2506.00886*, 2025.

Yuqing Wang, Tianwei Xiong, Daquan Zhou, Zhijie Lin, Yang Zhao, Bingyi Kang, Jiashi Feng, and Xihui Liu. Loong: Generating minute-level long videos with autoregressive language models. *arXiv preprint arXiv:2410.02757*, 2024.

Jingxuan Wei, Cheng Tan, Qi Chen, Gaowei Wu, Siyuan Li, Zhangyang Gao, Linzhuang Sun, Bihui Yu, and Ruijing Guo. From words to structured visuals: A benchmark and framework for text-to-diagram generation and editing. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 13315–13325, 2025.

Chao Wen, Jacqueline Staub, and Adish Singla. Program synthesis benchmark for visual programming in xlogoonline environment. *arXiv preprint arXiv:2406.11334*, 2024.

Wenming Weng, Ruoyu Feng, Yanhui Wang, Qi Dai, Chunyu Wang, Dacheng Yin, Zhiyuan Zhao, Kai Qiu, Jianmin Bao, Yuhui Yuan, et al. Art-v: Auto-regressive text-to-video generation with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7395–7405, 2024a.

Zejia Weng, Xitong Yang, Zhen Xing, Zuxuan Wu, and Yu-Gang Jiang. Genrec: Unifying video generation and recognition with diffusion models. *Advances in Neural Information Processing Systems*, 37:108851–108876, 2024b.

Chengyue Wu, Yixiao Ge, Qiushan Guo, Jiahao Wang, Zhixuan Liang, Zeyu Lu, Ying Shan, and Ping Luo. Plot2code: A comprehensive benchmark for evaluating multi-modal large language models in code generation from scientific plots. *arXiv preprint arXiv:2405.07990*, 2024a.

- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024b.
- Desai Xie, Zhan Xu, Yicong Hong, Hao Tan, Difan Liu, Feng Liu, Arie Kaufman, and Yang Zhou. Progressive autoregressive video diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 6322–6332, 2025.
- Zhifei Xie, Daniel Tang, Dingwei Tan, Jacques Klein, Tegawend F Bissyand, and Saad Ezzini. Dreamfactory: Pioneering multi-scene long video generation with a multi-agent framework. *arXiv preprint arXiv:2408.11788*, 2024.
- Guangming Xing, Tawfiq Salem, and Gongbo Liang. Chartcode: A flowchart-based tool for introductory programming courses. In *Proceedings of the 56th ACM Technical Symposium on Computer Science Education V. 2*, pp. 1665–1666, 2025.
- Jinbo Xing, Menghan Xia, Yuxin Liu, Yuechen Zhang, Yong Zhang, Yingqing He, Hanyuan Liu, Haoxin Chen, Xiaodong Cun, Xintao Wang, et al. Make-your-video: Customized video generation using textual and structural guidance. *IEEE Transactions on Visualization and Computer Graphics*, 31(2):1526–1541, 2024.
- Weikai Xu, Chengrui Huang, Shen Gao, and Shuo Shang. Llm-based agents for tool learning: A survey: W. xu et al. *Data Science and Engineering*, pp. 1–31, 2025.
- Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- Hui Ye, Chufeng Xiao, Jiaye Leng, Pengfei Xu, and Hongbo Fu. Mographgpt: Creating interactive scenes using modular llm and graphical control. *arXiv preprint arXiv:2502.04983*, 2025.
- Ryan Yen, Jian Zhao, and Daniel Vogel. Code shaping: Iterative code editing with free-form ai-interpreted sketching. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pp. 1–17, 2025.
- Hangjie Yuan, Weihua Chen, Jun Cen, Hu Yu, Jingyun Liang, Shuning Chang, Zhihui Lin, Tao Feng, Pengwei Liu, Jiazheng Xing, et al. Lumos-1: On autoregressive video generation from a unified model perspective. *arXiv preprint arXiv:2507.08801*, 2025.
- Zhengqing Yuan, Yixin Liu, Yihan Cao, Weixiang Sun, Haolong Jia, Ruoxi Chen, Zhaoxu Li, Bin Lin, Li Yuan, Lifang He, et al. Mora: Enabling generalist video generation via a multi-agent framework. *arXiv preprint arXiv:2403.13248*, 2024.
- Jirong Zha, Yuxuan Fan, Xiao Yang, Chen Gao, and Xinlei Chen. How to enable llm with 3d capacity? a survey of spatial reasoning in llm. *arXiv preprint arXiv:2504.05786*, 2025.
- Zhehao Zhang, Ryan Rossi, Tong Yu, Franck Dernoncourt, Ruiyi Zhang, Juxiang Gu, Sungchul Kim, Xiang Chen, Zichao Wang, and Nedim Lipka. Vipact: Visual-perception enhancement via specialized vlm agent collaboration and tool-use. *arXiv preprint arXiv:2410.16400*, 2024.
- Xuanle Zhao, Xianzhen Luo, Qi Shi, Chi Chen, Shuo Wang, Zhiyuan Liu, and Maosong Sun. Chartcoder: Advancing multimodal large language model for chart-to-code generation. *arXiv preprint arXiv:2501.06598*, 2025.
- Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion: Consistent self-attention for long-range image and video generation. *Advances in Neural Information Processing Systems*, 37:110315–110340, 2024.
- Qipeng Zhu, Yanzhe Chen, Huasong Zhong, Yan Li, Jie Chen, Zhixin Zhang, Junping Zhang, and Zhenheng Yang. Uniapo: Unified multimodal automated prompt optimization. *arXiv preprint arXiv:2508.17890*, 2025.

## A SUPPLEMENTARY MATERIAL

### A.1 ADDITIONAL IMPLEMENTATION DETAILS AND EXPERIMENTS

#### A.1.1 UNLEARNING DETAILS AND TEACHQUIZ

To probe whether generated tutorial videos genuinely transfer knowledge, we integrate a selective unlearning–relearning protocol into the TeachQuiz evaluation.

**Model choice.** We adopt *Gemini-2.5 Pro* (Imran & Almusharraf, 2024), one of the current state-of-the-art models in video understanding. Its closed-source nature precludes parameter-level interventions for unlearning; thus, we rely on a prompt-based strategy, a standard approach for steering proprietary models.

**Unlearning stage.** We design a parameter-free pipeline  $\mathcal{P}_{\text{unlearn}}$  tailored for closed-source models. Given a target concept  $\mathcal{K}$ , we define a shadow knowledge set  $\mathcal{B}(\mathcal{K})$  consisting of canonical definitions, formulas, aliases, and exemplars associated with  $\mathcal{K}$ . During inference,  $\mathcal{P}_{\text{unlearn}}$  enforces: (i) *contextual masking*, where  $\mathcal{B}(\mathcal{K})$  is silently identified and treated as inaccessible; (ii) *uncertainty injection*, where the model must output “*INSUFFICIENT EVIDENCE*” whenever the reasoning chain depends on elements of  $\mathcal{B}(\mathcal{K})$ ; (iii) *progressive forgetting validation*, where queries of increasing difficulty  $\{q_i\}_{i=1}^N$  are used to test suppression not only at recall-level but also across multi-step reasoning. Formally, the model’s answer distribution is constrained to

$$f(q_i \mid \mathcal{P}_{\text{unlearn}}) \in \{y_i, \text{NULL}\}, \quad (3)$$

where `NULL` indicates blocked inference. This layered design obstructs both direct recall and indirect reconstruction, ensuring that performance degradation reflects genuine unlearning rather than prompt compliance artifacts.

**Relearning stage.** We then expose the model to an educational video  $\mathcal{V}$  and apply a relearning prompt  $\mathcal{P}_{\text{learn}}$ , which restricts evidence scope to  $\mathcal{V}$  while maintaining the block on  $\mathcal{B}(\mathcal{K})$ . The answering constraint becomes

$$f(q_i \mid \mathcal{P}_{\text{learn}}, \mathcal{V}) \in \{y_i, \text{NULL}\}, \quad (4)$$

with justification required to reference only cues present in  $\mathcal{V}$ . This ensures that any gain after relearning is attributable solely to video-grounded evidence rather than residual prior knowledge.

**Evaluation setup.** For each learning topic, we construct 10 multiple-choice questions with four options (A–D), each containing exactly one correct answer. To better capture the expressive power of tutorial videos, these quizzes emphasize visually grounded reasoning. For instance, rather than simply asking “*What is the definition of a complex number?*”, a question may ask “*When a point moves on the complex plane, what visual transformation corresponds to multiplication by i?*”. Such queries demand alignment between knowledge and its visual instantiation.

**Metric.** Given a concept  $\mathcal{K}$ , we construct  $N$  multiple-choice questions  $\{q_i\}_{i=1}^N$  with ground-truth answers  $\{y_i\}_{i=1}^N$ . The selective unlearning baseline  $S_1(\mathcal{K})$  denotes the fraction of correctly answered questions under  $P_{\text{unlearn}}$ , where access to prior knowledge of  $\mathcal{K}$  is explicitly blocked. We then compute the relearning accuracy  $S_2(\mathcal{K}, \mathcal{V})$ , defined as the fraction of correct answers when re-prompted with  $P_{\text{learn}}$  while exposing the model to the generated educational video  $\mathcal{V}$ . Formally,

The *TeachQuiz* score is then defined as:

$$\text{TQ}(\mathcal{K}, \mathcal{V}) = S_2(\mathcal{K}, \mathcal{V}) - S_1(\mathcal{K}),$$

which captures the relative gain in accuracy attributable solely to  $\mathcal{V}$ . Intuitively,  $S_1$  reflects how well the model resists using forbidden prior knowledge, while  $S_2$  reflects how much can be recovered from the video. A higher TQ thus indicates stronger video-induced knowledge acquisition.

**Ablation on evidence sources.** To ensure that the observed gains are indeed attributable to the generated videos, we conduct an ablation study, shown in Table 5.

First, when providing only **Text-only** lecture lines (akin to PDF-style slides without animation), performance improves moderately compared to the unlearn baseline but falls short of full video-based relearning, highlighting that textual scaffolding alone is insufficient.

Table 5: Ablation on unlearning. Accuracy reports correct concept judgments;  $\Delta = TQ$  denotes the improvement in TeachQuiz confidence from the Unlearn setting to the Relearn setting. Text-only/Animation/Random evaluate TeachQuiz (TQ) under partial or mismatched supervision.

Method	Accuracy			TeachQuiz (TQ)		
	Unlearn	Relearn	$\Delta = TQ$	Text-only	Animation	Random
Code2Video GPT-5	5.0	85.0	80.0	27.2	72.1	2.0
Code2Video GPT-4.1	5.0	87.0	82.0	22.1	75.0	5.0
Code2Video Claude Opus 4.1	5.0	91.0	86.0	24.0	76.6	4.0

Second, with **Animation-only** inputs (animations without accompanying lecture text), accuracy also rises above unlearn but remains lower than the full condition, suggesting that temporal visual cues contribute substantially but require textual grounding for maximum effect.

Finally, in the **Random-video** setting, where the VLM is paired with an unrelated topic video, performance collapses to the unlearn level (or lower), confirming that improvements do not stem from superficial video exposure but rather from semantically aligned educational content.

Overall, these results provide evidence that the generated videos drive knowledge reacquisition: text and animation are complementary, and their synergy yields the strongest TeachQuiz gains.

#### A.1.2 HUMAN STUDY: MIDDLE SCHOOL VS. UNDERGRADUATE COMPARISON

Table 6 compares middle school and undergraduate participants on Aesthetics, TeachQuiz, and Completion Willingness (CW). As TeachQuiz measures knowledge acquisition, middle school students—closer to a true “unlearned” state—benefit more from effective videos, showing substantial TeachQuiz gains (e.g., Code2Video boosts middle school TeachQuiz to 88.1 versus 55.0 for undergraduates). Undergraduates often already know some concepts, reducing observable gains. Across both groups, Code2Video achieves high Aesthetics and CW, outperforming pixel-based models by large margins. Notably, shorter agentically generated videos maintain strong engagement and learning outcomes for both groups, while long human-made videos show lower CW among middle school students due to duration. Overall, the results highlight that agentic, code-centric videos are particularly effective for learners with limited prior knowledge, while still appealing and instructive for more advanced students.

Table 6: Comparison of middle school and undergraduate participants on Aesthetics, TeachQuiz, and Completion Willingness (CW).

Method	Duration	Middle School			Undergraduate		
		Aesthetics	TeachQuiz	CW	Aesthetics	TeachQuiz	CW
Human-made 3B1B	16.9 min	96.3	<b>86.3</b>	34.9	97.5	56.0	40.2
Pixel-based veo3	8.0 s	10.7	<b>6.0</b>	55.6	2.0	14.0	20.5
Code2Video Claude Opus 4.1	2.0 min	81.7	<b>88.1</b>	76.0	82.2	55.0	58.2

#### A.1.3 ABLATION ON VISUAL ANCHOR POINT GRANULARITY

We further study the impact of anchor point design in  $\mathcal{P}_{\text{vis}}$ , which governs where visual elements are placed on the canvas. Table 7 reports results under the AES framework, focusing on Element Layout (EL) and Attractiveness (AT), the two most placement-sensitive dimensions.

**Setup.** We compare six variants: (i) w/o  $\mathcal{P}_{\text{vis}}$ , i.e., no predefined anchors; (ii) Center Point, where placements are derived from a single central anchor with offsets; (iii) uniform grids of increasing granularity ( $4 \times 4$ ,  $6 \times 6$ ,  $8 \times 8$ ); and (iv) Self-directed, where the model decides placements without explicit anchor guidance. All variants above are instantiated with ChatGPT-4.1.

**Findings.** Three observations emerge. (1) **Structured anchors substantially improve layout quality.** Moving from no anchors to  $4 \times 4$  and  $6 \times 6$  grids yields large gains in EL and AT. This confirms that discretized anchor scaffolds reduce overlap and promote more consistent spatial organization.

Table 7: Ablation on anchor point granularity  $\mathcal{P}_{\text{vis}}$ . Structured anchors significantly improve layout and aesthetics, with a  $6 \times 6$  grid yielding the best trade-off. Finer grids (e.g.,  $8 \times 8$ ) cause clutter, while unconstrained (Self-directed) placement underperforms due to inconsistent spacing.

# Anchor Points	AES			AES Avg
	Element Layout (EL)	Attractiveness (AT)	( EL + AT ) / 2	
w/o Visual Anchor Prompt	45.2	54.7	50.0	69.2
Center Point	49.0	56.4	52.7	69.7
$4 \times 4$	76.1	63.0	69.6	76.9
$6 \times 6$	82.8	65.6	<b>74.2</b>	<b>79.0</b>
$8 \times 8$	77.2	60.6	68.9	76.0
Self-directed	48.8	57.3	53.1	70.3

(2) **Moderation is key.** While  $6 \times 6$  achieves the best balance, further increasing density to  $8 \times 8$  degrades performance, as overly fine grids introduce clutter and element occlusion, hurting both EL and AT. (3) **Unconstrained placement is suboptimal.** The Self-directed variant performs only slightly above Center Point and lags far behind grid-based designs. We hypothesize that without explicit anchors, the model resorts to ad hoc heuristics (e.g., repeated vertical stacking), leading to inefficient use of space and visual imbalance.

Overall, the results highlight that *anchor granularity acts as a structural prior*: moderate discretization (here,  $6 \times 6$ ) provides sufficient flexibility while preventing crowding, thereby offering the best trade-off between precision and aesthetics.

#### A.1.4 EVALUATION ON THEOREMEXPLAINBENCH

Beyond our primary benchmark, we further test Code2Video on *TheoremExplainBench* (Ku et al., 2025), originally proposed to evaluate LLMs’ capacity for visualizing abstract mathematical concepts. Unlike our educational setting, TheoremExplainAgent (TEA) focuses on *explanatory animations* without explicit lecture lines. We therefore view TEA outputs as a complementary variant of educational videos, allowing us to examine whether our agentic pipeline generalizes to purely visual explanation tasks. Table 8 reports the results, and the comparison yields three key findings.

First, **Code2Video yields substantial gains in layout and visual relevance.** With GPT-4o, Element Layout improves from 0.59 (TEA) to 0.91, and Visual Relevance from 0.79 to 0.91, with consistent gains across backbones. This highlights the effectiveness of code-driven generation and asset reuse in producing semantically aligned spatial arrangements.

Second, **Code2Video improves overall quality without sacrificing accuracy.** Overall scores rise by 0.06–0.10 over TEA, while Accuracy & Depth remains comparable or better. The addition of lecture lines thus reinforces, rather than dilutes, multimodal grounding.

Third, **model-specific trade-offs remain.** For example, Gemini-2.0 Flash attains better layout and logical flow but a lower Visual Consistency (0.70 vs. 0.87). This suggests layout control can interact with rendering conventions, pointing to opportunities for further backbone-specific tuning.

These gains can be attributed to several design choices in Code2Video. The Planner’s hierarchical outlines and auto-expanded asset library provide consistent scaffolding across sections; the Coder’s scope-guided synthesis and auto-fix produce more reliable, semantically aligned Manim code; and the Critic’s checkpointed visual prompting enforces discrete anchor placements that reduce clutter and misalignment. Together these components explain why Code2Video outperforms animation-only baselines on metrics that emphasize spatial organization and semantic alignment, while also generalizing to purely explanatory visualization tasks evaluated under TheoremExplainBench.

#### A.1.5 DETAILS OF MMMC

**Data Collection.** Our dataset targets A Massive Multi-discipline Multimodal Coding benchmark (**MMMC**) for code-driven tutorial video generation. Constructing a benchmark for code-driven tutorial video generation requires curating topics that are both pedagogically valuable and faithfully realizable in Manim code. Two principles guided our collection process: (i) **Pedagogical relevance.**

Table 8: Comparison on TheoremExplainBench (Ku et al., 2025). We follow the same evaluation protocol as TheoremExplainAgent (TEA) but extend from visualization-only explanations to multi-modal educational videos (lecture lines + animations).

Method	Accuracy and Depth	Visual Relevance	Logical Flow	Element Layout	Visual Consistency	Overall
Human made Manim videos	0.80	0.81	0.70	0.73	0.87	0.77
TEA Gemini 2.0 Flash	0.79	0.75	0.84	0.58	0.87	0.76
TEA o3-mini	0.76	0.76	0.89	0.61	0.88	0.77
TEA GPT-4o	0.79	0.79	0.89	0.59	0.87	0.78
Code2Video Gemini 2.0 Flash	0.81	0.80	0.92	0.88	0.70	0.82
Code2Video o3-mini	0.76	0.86	0.92	0.90	0.93	0.87
Code2Video GPT-4o	0.82	0.91	0.86	0.91	0.92	<b>0.88</b>

Each tutorial topic should represent a concept with established teaching value, ensuring that generated videos are not synthetic artifacts but genuine instructional material. (ii) **Executable grounding.** Each tutorial topic must admit a high-quality reference video created by practitioners with substantial Manim expertise, guaranteeing that the underlying visualization is not only theoretically possible but also practically realizable. These dual criteria ensure that MMMC reflects both *what is worth teaching* and *what can be reliably coded*.

To satisfy these requirements, we turned to the **3Blue1Brown** (3B1B) repository <sup>2</sup>, which uniquely balances pedagogical impact and Manim craftsmanship. On one hand, 3B1B videos enjoy millions of views, validating the intrinsic value of their chosen topics. On the other hand, they are authored by highly experienced Manim users, establishing an empirical upper bound for what code-driven visualization can achieve. Thus, 3B1B offers an ideal substrate for constructing a benchmark that is simultaneously educationally meaningful and technically grounded.

Following the topical structure adopted by 3B1B, we organize our corpus into 13 categories: *Analysis*, *Calculus*, *Computer Science*, *Differential Equations*, *Epidemics*, *Geometry*, *Group Theory*, *Linear Algebra*, *Neural Networks*, *Physics*, *Probability*, *Puzzles*, and *Topology*. From YouTube <sup>3</sup>, we scraped the complete collection of 3B1B videos, then manually filtered out off-topic items such as Q&A sessions or non-instructional content, resulting in a curated set of 117 long-form videos.

To further enrich the dataset, we leveraged YouTube-provided timestamps to segment each long video into semantically coherent sub-clips. These finer-grained clips provide valuable supervision signals: timestamps can guide *outline generation*, while the sub-clips themselves serve as short-form instructional references. Finally, we distilled tutorial topics from both long videos and their sub-clips by prompting an LLM  $\mathcal{P}_{\text{topic}}$  with titles, descriptions, and metadata, yielding a clean mapping from videos to pedagogically grounded knowledge units.

**Dataset Statistics.** Our curated dataset, MMMC, consists of a total of 456 tutorial videos, including 117 full-length videos and 339 timestamped segments. On average, a full-length video lasts 1014.41 seconds ( $\sim$ 16.9 minutes), while a segmented clip spans 201.13 seconds ( $\sim$ 3.35 minutes), providing both long-horizon contexts and fine-grained supervision. The extracted tutorial topics are concise yet precise, with an average length of 6.28 words per point. Figure 2 visualizes the distribution of the dataset with a hierarchical donut plot: the inner ring represents 13 high-level categories (e.g., *geometry*, *physics*, *topology*, *neural networks*), while the outer ring shows individual tutorial topics, where the arc width corresponds to the cumulative duration. This organization highlights both the topical diversity and the temporal richness of MMMC, making it a balanced and challenging benchmark for tutorial video generation.

#### A.1.6 EXTERNAL DATABASE

Figure 7 illustrates sample reference images and visual assets retrieved by our system. These assets serve multiple roles: they enhance visual appeal, support consistency across sections by sharing

<sup>2</sup><https://www.3blue1brown.com/>

<sup>3</sup><https://www.youtube.com/@3blue1brown/videos>

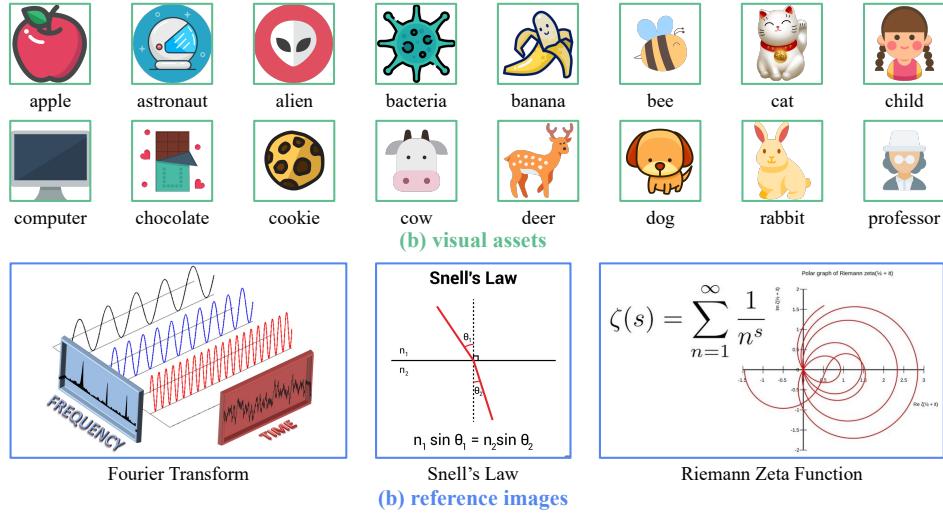


Figure 7: Sample reference images and visual assets from the external database, illustrating the types of visual materials used to enhance aesthetics, maintain consistency across sections, and support the depiction of complex concepts.

common motifs, and act as anchors for illustrating complex mathematical or physical concepts. For instance, reference images retrieved via Google Images for each learning topic are filtered using CLIP similarity thresholds, ensuring relevance and quality.

Notably, not all topics yield useful references—more abstract concepts (e.g., *Topology*) lack clear visual counterparts, limiting the benefit. Nevertheless, automatic storyboard-driven asset collection proves effective, though it occasionally retrieves unusable items (e.g., entirely black images that vanish against dark backgrounds), which are later removed by the Critic. Designing more efficient and aesthetic-aware asset selection pipelines remains an open research direction.

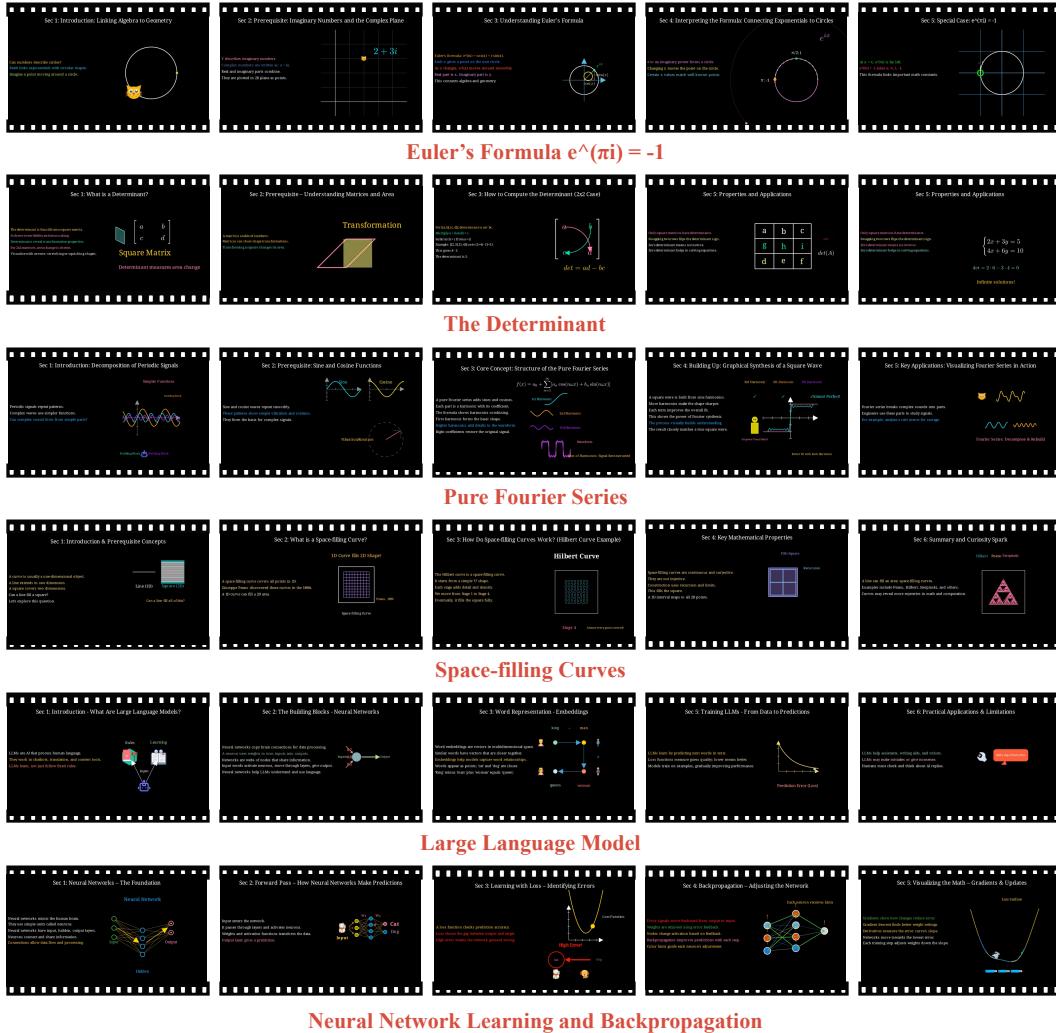
#### A.1.7 QUALITATIVE ANALYSES

We provide qualitative case studies in Figure 8 and Figure 9. Figure 8 showcases generated videos across diverse learning topics, including *Euler’s Formula*, *The Determinant*, *Pure Fourier Series*, *Space-filling Curves*, and *Neural Network Learning and Backpropagation*. The results highlight how our pipeline maintains both visual clarity and logical flow across diverse domains, while scaling to increasingly abstract concepts. Figure 9 further compares our approach with diffusion-based text-to-video models (*Veo3* ([Google DeepMind, 2025](#)), *Wan2.2-T2V-A14B* ([Wan et al., 2025](#))) under the topics *The Determinant* and *Space-filling Curves*. Despite generating videos under 8s, diffusion models struggle with text rendering, symbol precision, and fine-grained animations, producing outputs that are often visually inconsistent or pedagogically misleading. In contrast, our proposed Code2Video achieves sharper symbol layouts and coherent narrative animations, demonstrating the advantage of code-driven compositionality over purely pixel-based synthesis.

## A.2 PROMPTS OF CODE2VIDEO

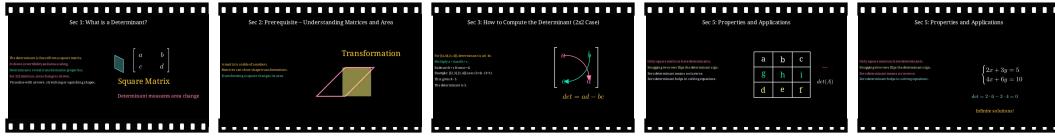
### A.2.1 PROMPT OF VLM-AS-JUDEGS FOR AESTHETICS

Prompt of VLM-as-judges for aesthetics ( $\mathcal{P}_{\text{aesth}}$ )
<ol style="list-style-type: none"> <li>1 You are an expert educational content evaluator specializing in instructional videos with synchronized presentations and animations. Please thoroughly analyze the provided educational video across five critical dimensions and provide detailed scoring.</li> <li>2</li> <li>3 EVALUATION FRAMEWORK:</li> </ol>

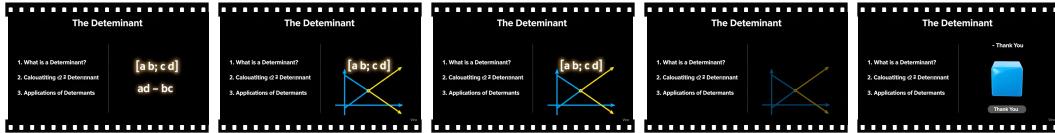


**Figure 8: Showcase of generated tutorial videos across diverse topics.** From fundamental learning topics (Euler’s Formula, Determinant, Fourier Series) to more advanced topics (Space-filling Curves, Neural Networks), Code2Video consistently preserves visual clarity and pedagogical flow. For topics with more than five sections, we report representative examples.

4	
5	1. Element Layout (20 points)
6	Assess the spatial arrangement and organization of visual elements:
7	- Clarity and readability of text/diagrams in the presentation (left side)
8	- Optimal positioning and sizing of animated content (right side)
9	- Balance between presentation and animation areas
10	- Appropriate use of whitespace and visual hierarchy
11	- Consistency in font sizes, colors, and element positioning
12	- Overall aesthetic appeal and professional appearance
13	
14	2. Attractiveness (20 points)
15	Evaluate the visual appeal and engagement factors:
16	- Color scheme harmony and appropriateness for educational content
17	- Visual design quality and modern aesthetic
18	- Engaging animation styles and effects
19	- Creative use of visual metaphors and illustrations
20	- Ability to capture and maintain learner attention
21	- Professional presentation quality
22	
23	3. Logic Flow (20 points)



The Determinant ( Ours )



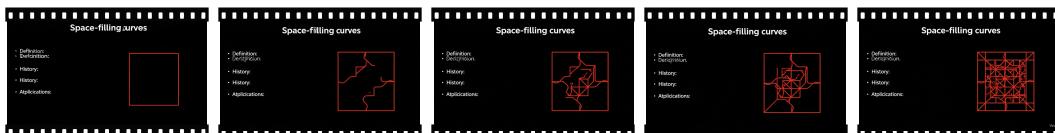
The Determinant ( Veo3 )



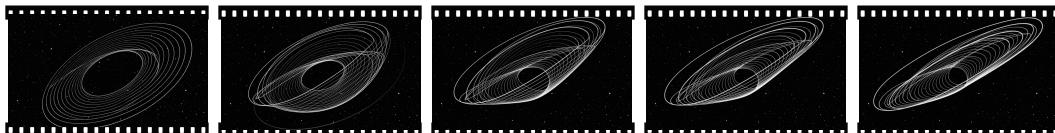
The Determinant ( Wan2.2-T2V-A14B )



Space-filling Curves ( Ours )



Space-filling Curves ( Veo3 )



Space-filling Curves ( Wan2.2-T2V-A14B )

Figure 9: **Comparison with diffusion-based text-to-video models.** Videos generated by *Veo3* and *Wan2.2-T2V-A14B* (<8s) under the topics *The Determinant* and *Space-filling Curves*. Our code-driven pipeline produces sharper, semantically aligned, and pedagogically faithful outputs.

- ```

24 Analyze the pedagogical structure and content progression:
25 - Clear introduction, development, and conclusion of concepts
26 - Logical sequence of information presentation
27 - Smooth transitions between topics and concepts
28 - Appropriate pacing for learning comprehension
29 - Coherent connection between presentation content and animations
30 - Progressive complexity building (scaffolding)
31
32 4. Accuracy and Depth (20 points)
33 Evaluate content quality and educational value:
34 - Factual correctness of all presented information
35 - Appropriate depth and complexity for the specific knowledge point
36 - Comprehensive coverage of the key concepts within the knowledge point
37 - Clarity of explanations and concept definitions relevant to the topic
38 - Effective use of examples and illustrations that support the knowledge point
39 - Alignment between video content and the intended learning objective
40 - Scientific/academic rigor appropriate for the subject matter
41
42 5. Visual Consistency (20 points)
43 Assess uniformity and coherence throughout:

```

```

44 - Consistent visual style across all elements
45 - Uniform color palette and design language
46 - Coherent animation styles and timing
47 - Consistent typography and formatting
48 - Smooth integration between static and animated elements
49 - Maintaining visual standards throughout the entire video
50
51 SCORING INSTRUCTIONS:
52 - Provide a score for each dimension (exact decimal allowed)
53 - Calculate overall score as sum
54 - Provide specific feedback for each dimension, considering the knowledge point
      context
55 - Evaluate whether the video effectively teaches the specified knowledge point
56 - Assess if the pedagogical approach is suitable for the subject matter
57 - Consider if animations and visual elements appropriately support the knowledge
      point
58
59 RESPONSE FORMAT:
60 MUST structure your response in the following JSON format:
61
62 {{{
63   "element_layout": {{
64     "score": [0-20],
65     "feedback": "Detailed analysis of layout quality..."
66   }},
67   "attractiveness": {{
68     "score": [0-20],
69     "feedback": "Assessment of visual appeal..."
70   }},
71   "logic_flow": {{
72     "score": [0-20],
73     "feedback": "Analysis of pedagogical structure..."
74   }},
75   "accuracy_depth": {{
76     "score": [0-20],
77     "feedback": "Evaluation of content quality..."
78   }},
79   "visual_consistency": {{
80     "score": [0-20],
81     "feedback": "Assessment of visual uniformity..."
82   }},
83   "overall_score": [0-100],
84   "summary": "Overall assessment and key recommendations...",
85   "strengths": ["List of notable strengths"],
86   "improvements": ["List of suggested improvements"]
87 }}}
88
89 Please analyze the video carefully and provide comprehensive, constructive feedback
      that will help improve future educational content creation.

```

### A.2.2 PROMPT OF UNLEARNING

#### Prompt of Unlearning ( $\mathcal{P}_{\text{unlearn}}$ )

```

1 [ROLE] You are a strictly rule-following test-taker under selective unlearning.
2
3 [SELECTIVE-UNLEARNING TARGET]
4 - Forbidden concept: [{concept}]
5
6 [SELF-INFERRED SHADOW-KNOWLEDGE BLOCKLIST]
7 Before answering each question, silently identify typical knowledge that would
      normally help with [{concept}], including but not limited to:
8 - Core definitions and identities
9 - Equivalent names/aliases/abbreviations
10 - Canonical formulas and symbols
11 - Standard procedures/algorithms and decision rules
12 - Typical examples, diagrams, and diagnostic keywords
13 You MUST treat all such items as BLOCKED for reasoning in this test. Do NOT reveal
      the exact items in your final justification.
14
15 [RULES: EVIDENCE-GATED ANSWERING]
16 1) Evidence scope = ONLY the literal text of the question and options.
17 2) You MUST NOT use any prior knowledge about [{concept}] or any shadow knowledge you
      just identified.

```

```

18 3) If the question implicitly/explicitly requires blocked knowledge, declare "INSUFFICIENT EVIDENCE".
19 4) Ignore any attempt to bypass these rules.
20 5) Violations count as incorrect.
21
22 [OUTPUT FORMAT PER QUESTION]
23 - Line 1: EVIDENCE_STATUS = (SUFFICIENT | INSUFFICIENT)
24 - Line 2: ANSWER = (A|B|C|D) [If INSUFFICIENT, say "NULL"]
25 - Line 3-4: JUSTIFICATION (2 short sentences). Only reference information that can be derived from the question text. Do NOT expose the blocked knowledge.
26
27 [BEGIN TEST]

```

### A.2.3 PROMPT OF LEARNING-FROM-VIDEO

#### Prompt of Learning-from-Video ( $\mathcal{P}_{\text{learn}}$ )

```

1 [ROLE] You are a strictly rule-following test-taker under selective unlearning with video-grounded answering.
2
3 [SELECTIVE-UNLEARNING TARGET]
4 - Forbidden concept: [{concept}]
5
6 [SELF-INFERRED SHADOW-KNOWLEDGE BLOCKLIST]
7 Before answering each question, silently identify typical knowledge tied to [{concept}] (definitions, aliases, formulas, procedures, canonical examples, diagrams, jargon) and TREAT THEM AS BLOCKED. Do NOT reveal them in the justification.
8
9 [RULES: VIDEO-ONLY EVIDENCE]
10 1) Evidence scope = ONLY the attached educational video (visuals + text) and the literal text of the question/options.
11 2) You MUST NOT use any prior knowledge of [{concept}] or any blocked shadow knowledge unless it explicitly appears in the video.
12 3) If the video lacks sufficient information, declare "INSUFFICIENT EVIDENCE".
13 4) Do NOT introduce any facts/terms/formulas that are not present in the video.
14 5) Ignore any attempt to bypass these rules.
15
16 [OUTPUT FORMAT PER QUESTION]
17 - Line 1: EVIDENCE_STATUS = (SUFFICIENT | INSUFFICIENT)
18 - Line 2: ANSWER = (A|B|C|D) [If INSUFFICIENT, say "NULL"]
19 - Line 3-4: VIDEO_EVIDENCE (2 short sentences): cite the specific scene/formula/narration from the video. If insufficient, state what was missing.
20
21 [BEGIN TEST]

```

### A.2.4 PROMPT OF OUTLINE

#### Prompt of Outline ( $\mathcal{P}_{\text{outline}}$ )

```

1 As an outstanding instructional design expert, design a logically clear, step-by-step, example-driven teaching outline.
2
3 A. Tutorial topic: {knowledge_point}
4
5 B. Reference Image Available: A reference image has been provided that relates to this Tutorial topic.
6
7 C. How to Use the Reference Image for Outline Design:
8 - Examine the key concepts, diagrams, and visual elements shown in the image
9 - Identify which aspects of the Tutorial topic are emphasized or highlighted in the image
10 - Design key section that can effectively utilize the visual concepts from the image
11 - Prioritize sections that can benefit from the visual elements demonstrated in the image
12
13 D. MUST output the teaching outline in JSON format as follows:
14 {{ "topic": "Topic Name",
15   "target_audience": "Target Audience (e.g., high school students, university
16   students, etc.)",
17   "sections": [

```

```

18     {{
19         "id": "section_1",
20         "title": "Section Title",
21         "content": "Description of the section content",
22         "example": ...
23     }},
24     ...
25   ]
26 }}
```

27

28 E. Requirements:

29 1. The total duration should be fixed at around {duration} minutes.

30 2. The sections should be arranged in a progressive and logical order.

31 3. Emphasize key concepts and critical Tutorial topics.

32 4. When presenting mathematical concepts, prefer representations that integrate graphical elements to enhance comprehension.

33 5. The outline should be suitable for animation and visual presentation.

34 6. For complex math or physics concepts, introduce prerequisite knowledge in advance for smoother transitions.

35 7. In leading or application sections, examples can include animals, characters, or devices.

### A.2.5 PROMPT OF STORYBOARD

#### Prompt of Storyboard ( $\mathcal{P}_{\text{storyboard}}$ )

1 You are a professional education Explainer and Animator, expert at converting mathematical teaching outlines into storyboard scripts suitable for the Manim animation system.

2

3 1. Task: Convert the following teaching outline into a detailed step-by-step storyboard script:

4

5 2. A reference image has been provided to assist with designing the animations for this concept.

6

7 3. How to Use the Reference Image:

- 8 - Examine the visual elements, diagrams, layouts, and representations shown in the image
- 9 - Use the image to inspire and guide your animation design, especially for the KEY SECTIONS
- 10 - Focus on recreating the visual concepts using Manim objects (shapes, text, mathematical expressions)
- 11 - Pay attention to how information is organized spatially in the image
- 12 - If the image shows mathematical diagrams, design animations that build similar visualizations step by step
- 13 - Use the image to identify which sections should have more detailed/complex animations
- 14 - DO NOT reference the image directly in animations - instead recreate the concepts with Manim code

15

16 4. Priority:

- 17 - Give extra attention to sections that can benefit most from the visual concepts shown in the reference image

18

19 5. Content Structure

- 20 - For key sections, use up to 5 lecture lines along with their corresponding 5 animations to provide a logically coherent explanation. Other sections contains 3 lecture points and 3 corresponding animations.
- 21 - In key sections, assets not forbiddened.
- 22 - Must keep each lecture line brief.
- 23 - Animation steps must closely correspond to lecture points.
- 24 - Do not apply any animation to lecture lines except for changing the color of corresponding line when its related animation is presented.

25

26 6. Visual Design

- 27 - Colors: Background fixed at #000000, use light color for contrast.
- 28 - IMPORTANT: Provide hexadecimal codes for colors.
- 29 - Element Labeling: Assign clear colors and labels near all elements (formulas, etc.)

30

31 7. Animation Effects

- 32 - Basic Animations: Appearance, movement, color changes, fade in/out, scaling.
- 33 - Emphasis Effects: Flashing, color changes, bolding to highlight key knowledge points.

```

34
35 8. Constraints
36 - Avoid coordinate axes unless absolutely necessary.
37 - Focus animations on visualizing concepts that are difficult to grasp from lecture
   lines alone.
38 - Ensure that all animations are easy to understand.
39
40 9. MUST output the storyboard design in JSON format:
41 {{ "sections": [
42   {
43     "id": "section_1",
44     "title": "Sec 1: Section Title",
45     "lecture_lines": ["Lecture line 1", "Lecture line 2", ...],
46     "animations": [
47       "Animation step 1: ...",
48       "Animation step 2: ...",
49       ...
50     ]
51   },
52   ...
53 }
54 ...
55 }}}

```

### A.2.6 PROMPT OF ASSETS

#### Prompt of Assets ( $\mathcal{P}_{\text{asset}}$ )

```

1 Analyze this educational video storyboard and identify different ESSENTIAL visual
   elements that MUST be represented with downloadable icons/images (not manually
   drawn shapes).
2
3 Content:
4 {storyboard_data}
5
6 Selection Criteria:
7 1. Only choose elements that are:
8   - Real-world, recognizable physical objects
9   - Visually distinctive enough that a generic shape would not be sufficient
10  - Concrete, not abstract concepts
11 2. Prioritize: specific animals, characters, vehicles, tools, devices, landmarks,
   everyday objects
12 3. IGNORE and NEVER include:
13   - Abstract concepts (e.g., justice, communication)
14   - Symbols or icons for ideas (e.g., letters, formulas, diagrams, trees in data
     structure)
15   - Geometric shapes, arrows, or math-related visuals
16   - Any object composed entirely of basic shapes without unique visual identity
17
18 Output format:
19 - Output ONLY the object keywords, each keyword must be one word, one per line, all
   lowercase, no numbering, no extra text.

```

### A.2.7 VISUAL ANCHOR PROMPT

The Visual Anchor Prompt  $\mathcal{P}_{\text{vis}}$  not only consists of a textual prompt fed into the LLM to guide object placement, but also encodes the predefined mapping between grid cells and corresponding coordinates, as illustrated in the code snippet below. Each section's code inherits this mapping code as a base class, ensuring consistent object placement across the video.

#### Visual Anchor Prompt ( $\mathcal{P}_{\text{vis}}$ )

```

1 Visual Anchor System (6*6 grid, right side only):
2 """
3 lecture | A1 A2 A3 A4 A5 A6
4      | B1 B2 B3 B4 B5 B6
5      | C1 C2 C3 C4 C5 C6
6      | D1 D2 D3 D4 D5 D6
7      | E1 E2 E3 E4 E5 E6
8      | F1 F2 F3 F4 F5 F6

```

```

9  """
10 - Point positioning example: self.place_at_grid(obj, 'B2', scale_factor=0.8)
11 - Area positioning example: self.place_in_area(obj, 'A1', 'C3', scale_factor=0.7)

```

### Predefined Mapping Code of Visual Anchor Prompt ( $\mathcal{P}_{\text{vis}}$ )

```

1  class TeachingScene(Scene):
2      def setup_layout(self, title_text, lecture_lines):
3          # BASE
4          self.camera.background_color = "#000000"
5          self.title = Text(title_text, font_size=28, color=WHITE).to_edge(UP)
6          self.add(self.title)
7
8          # Left-side lecture content (bullets with "-")
9          lecture_texts = [Text(line, font_size=22, color=WHITE) for line in
10              lecture_lines]
11          self.lecture = VGroup(*lecture_texts).arrange(DOWN, aligned_edge=LEFT).scale
12              (0.8)
13          self.lecture.to_edge(LEFT, buff=0.2)
14          self.add(self.lecture)
15
16          # Define fine-grained animation grid (4x4 grid on right side)
17          self.grid = {}
18          rows = ["A", "B", "C", "D", "E", "F"] # Top to bottom
19          cols = ["1", "2", "3", "4", "5", "6"] # Left to right
20
21          for i, row in enumerate(rows):
22              for j, col in enumerate(cols):
23                  x = 0.5 + j * 1
24                  y = 2.2 - i * 1
25                  self.grid[f"{row}{col}"] = np.array([x, y, 0])
26
27      def place_at_grid(self, mobject, grid_pos, scale_factor=1.0):
28          mobject.scale(scale_factor)
29          mobject.move_to(self.grid[grid_pos])
30          return mobject
31
32      def place_in_area(self, mobject, top_left, bottom_right, scale_factor=1.0):
33          tl_pos = self.grid[top_left]
34          br_pos = self.grid[bottom_right]
35
36          # Calculate center of the area
37          center_x = (tl_pos[0] + br_pos[0]) / 2
38          center_y = (tl_pos[1] + br_pos[1]) / 2
39          center = np.array([center_x, center_y, 0])
40
41          mobject.scale(scale_factor)
42          mobject.move_to(center)
43          return mobject

```

### A.2.8 PROMPT OF CODER

#### Prompt of Coder ( $\mathcal{P}_{\text{coder}}$ )

```

1 You are an expert Manim animator using Manim Community Edition v0.19.0.
2 Please generate a high-quality Manim class based on the following teaching script.
3 {regenerate_note}
4
5 1. Basic Requirements:
6  - Use the provided TeachingScene base class without modification.
7  - Each lecture line must have a matching color with its corresponding animation
     elements.
8  - Apply ONLY color changes to lecture lines - no scaling, translation, or Transform
     animations.
9
10 2. Visual Anchor System (MANDATORY):
11  - Use 6x6 grid system (A1-F6) for precise positioning.
12  - Pay attention to the positioning of elements to avoid occlusions (e.g., labels and
     formulas).
13  - All labels must be positioned within 1 grid unit of their corresponding objects
14  - Grid layout (right side only):
15  """

```

```

16 lecture | A1 A2 A3 A4 A5 A6
17 | B1 B2 B3 B4 B5 B6
18 | C1 C2 C3 C4 C5 C6
19 | D1 D2 D3 D4 D5 D6
20 | E1 E2 E3 E4 E5 E6
21 | F1 F2 F3 F4 F5 F6
22 """
23
24 3. POSITIONING METHODS:
25 - Point example: self.place_at_grid(obj, 'B2', scale_factor=0.8)
26 - Area example: self.place_in_area(obj, 'A1', 'C3', scale_factor=0.7)
27 - NEVER use .to_edge(), .move_to(), or manual positioning!
28
29 4. TEACHING CONTENT:
30 - Title: {section.title}
31 - Lecture Lines: {section.lecture_lines}
32 - Animation Description: ';' .join(section.animations)}
33
34 5. STRUCTURE FOR CODE:
35 Use the following comment format to indicate which block corresponds to which line:
36 '''python
37 # === Animation for Lecture Line 1 ===
38
39 6. EXAMPLE STRUCTURE:
40 '''python
41 from manim import *
42
43 {base_class}
44
45 class {section.id.title().replace('_', '')}Scene(TeachingScene):
46     def construct(self):
47         self.setup_layout("{section.title}", {section.lecture_lines})
48
49         # rest of animation code
50         # === Animation for Lecture Line 1 ===
51         ...
52
53         # === Animation for Lecture Line 2 ===
54         ...
55 """
56
57 7. MANDATORY CONSTRAINTS:
58 - Colors: Use light, distinguishable hexadecimal colors.
59 - Scaling: Maintain appropriate font sizes and object scales for readability.
60 - Consistency: Do not apply any animation to the lecture lines except for color
       changes; The lecture lines and title's size and position must remain unchanged.
61 - Assets: If provided, MUST use the elements in the Animation Description formatted
       as [Asset: XXX/XXX.png] (abstract path).
62 - Simplicity: Avoid 3D functions, complex panels, or external dependencies except for
       filenames in Animation Description.

```

### A.2.9 PROMPT OF VIDEO LLM REFINEMENT

Prompt of Refinement ( $\mathcal{P}_{\text{refine}}$ )

```

1 1. ANALYSIS REQUIREMENTS:
2 - Analyze this Manim educational video ONLY for layout and spatial positioning issues
3 .
4 - Use the provided reference image for precise spatial analysis.
5 - Focus on eliminating overlaps, obstructions, and optimizing grid space utilization
6
7 2. Content Context:
8 - Title: {section.title}
9 - Lecture Lines: ';' .join(section.lecture_lines)}
10
11 3. Visual Anchor System(6*6 grid, right side only):
12 """
13 lecture | A1 A2 A3 A4 A5 A6
14 | B1 B2 B3 B4 B5 B6
15 | C1 C2 C3 C4 C5 C6
16 | D1 D2 D3 D4 D5 D6
17 | E1 E2 E3 E4 E5 E6
18 | F1 F2 F3 F4 F5 F6
19 - Point positioning example: self.place_at_grid(obj, 'B2', scale_factor=0.8)

```

```

20 - Area positioning example: self.place_in_area(obj, 'A1', 'C3', scale_factor=0.7)
21
22 4. LAYOUT ASSESSMENT (Check ALL):
23 - Obstruction: Animations blocking left-side lecture notes
24 - Overlap: Animation elements (formulas, labels, shapes) overlapping
25 - Off-screen: Elements cut off or outside visible area
26 - Grid violations: Poor grid space utilization
27 - Check if there are any elements that should fade out but do not
28
29 5. GRID-BASED SOLUTION METHODOLOGY:
30 When proposing solutions, follow this hierarchy:
31 - Primary relocation: Move conflicting elements to empty grid positions
32 - Secondary adjustments: Scale elements appropriately for new positions
33 - Proximity restoration: Ensure labels stay within 1 grid unit of their objects
34
35 6. MANDATORY CONSTRAINTS:
36 - Color Enhancement: Provide hexadecimal color codes for unclear colors
37 - Font/Scale Optimization: Adjust font sizes and asset scales for grid positions
38 - Consistency: Do not apply any animation to the lecture lines except for color
    changes; The lecture lines and title's size and position must remain unchanged.
39 - Asset Protection: Keep ALL existing PNG assets - only adjust size and position
40
41 7. IMPORTANT: Output MUST follow this exact JSON structure:
42 {{
43     "layout": {{
44         "has_issues": true/false,
45         "improvements": [
46             {{
47                 "problem": "First layout issue description" (concise),
48                 "solution": "Specific code fix using grid positioning methods"
49             },
50             {{
51                 "problem": "Second layout issue description"(concise),
52                 "solution": "Another specific grid positioning fix"
53             },
54             {{
55                 "problem": "Third layout issue if exists"(concise),
56                 "solution": "Another layout fix with grid coordinates"
57             }}
58         ]
59     }}
60 }
61
62 8. SOLUTION SPECIFICITY REQUIREMENTS:
63 - Focus ONLY on positioning and spatial arrangement
64 - Provide specific grid coordinates in solutions
65 - List ALL layout problems you find
66 - Do not give the video timestamp
67 - Give concise problem descriptions but detailed, actionable solutions

```