



CODE2VIDEO: A CODE-CENTRIC PARADIGM FOR EDUCATIONAL VIDEO GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

While recent generative models have advanced pixel-space video synthesis, they remain limited in producing professional educational videos, which require 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 code while incorporating scope-guided auto-fix to enhance efficiency; and (iii) *Critic*, which uses a vision-language model (VLM) with visual anchor prompts to refine spatial layout and ensure clarity. For systematic evaluation, we construct **MMMC**, a benchmark of professionally produced, discipline-specific educational videos. We evaluate Code2Video 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 an ‘unlearned’ VLM 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.

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 high visual quality and short-form fidelity. Yet these models struggle with tasks that require long-form reasoning or complex multi-entity interactions (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 designed to teach subject-specific knowledge face unique challenges. 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; Fencel, 2010) 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.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

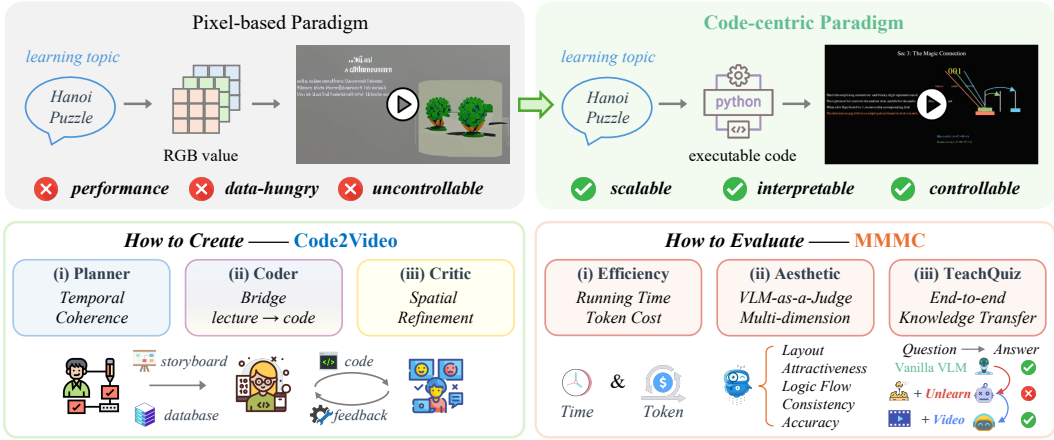


Figure 1: Overview of **Code2Video**. A code-centric paradigm for educational video generation, where Planner ensures temporal flow, Coder bridges instructions to executable animations, and Critic refines spatial layout. Evaluation is performed on **MMMC** with multi-dimensional metrics.

We are motivated by the unique advantages of code 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. Our framework decomposes the task into three collaborative agents: the *Planner* sequences concepts, examples, and recaps into a coherent lecture flow; the *Coder* translates structured instructions into executable Manim code, producing 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 areas (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 first unlearns a target concept from a VLM, and then measures how effectively the generated video restores that knowledge. Our experiments show several key findings: pixel-based models struggle with fine details and coherence, while a direct code-centric generation baseline improves TeachQuiz by 30%. Our full pipeline further achieves a stable 40% gain. In human studies based on TeachQuiz scores, videos from our pipeline even outperform professional human-made tutorials, demonstrating the power of our code-centric, agent-based approach.

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 propose a modular agent framework with three key components: (i) **The Planner** expands an external database for reference, enabling parallel yet consistent storyboard; (ii) **The Coder** ensures executable code via automatic debugging and scope-guided repair; (iii) **The Critic** refines spatial layout and clarity using visual anchor prompts.
- **A New Benchmark with Well-designed Evaluation Protocol.** 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, which makes it difficult to achieve the precise layout and symbolic alignment critical for educational videos. (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 *underexplored* (Ku et al., 2025; Wu et al., 2024b). (iii) Building on this paradigm, we propose a **code-centric animation framework** for educational video synthesis. Using executable code as the medium for generation enables symbolic layout, temporally structured exposition, and deterministic reproducibility—capabilities that are difficult to achieve with pixel-level diffusion.

2.2 CODING AGENTS

Recent advances in LLM-based tool use demonstrate that agents can autonomously call APIs, retrieve information, and verify outputs. This capability enables 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

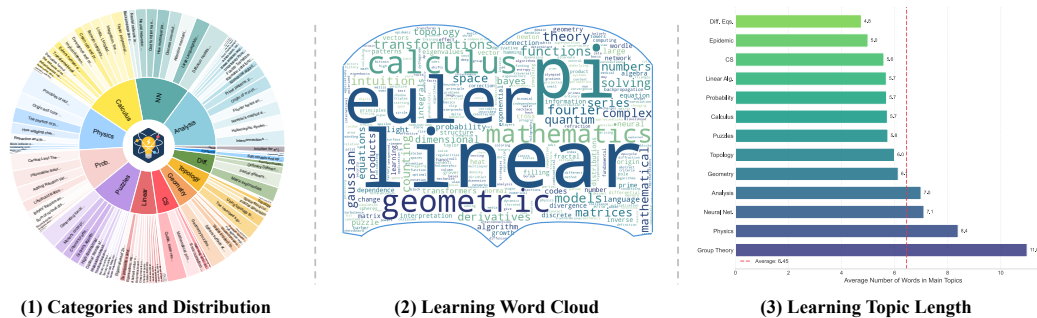


Figure 2: **MMMC overview**. (1) Left: Distribution of 13 subject areas with exemplar learning topics; the ring width represents video duration. Please refer to Fig. 7 for a clearer diagram. (2) Middle: **Word cloud of learning topics**. (3) Right: Average learning topic length subject per area.

162 3.1 TASK FORMULATION
163

164 The task of code-centric educational video generation maps a learning query to executable
165 *Manim* (Manim Community Dev, 2025) code whose rendering yields a tutorial video. The chal-
166 lenge lies in multi-step reasoning, precise temporal sequencing, and spatial coherence, where minor
167 syntax errors can prevent successful rendering. We adopt *Manim* for its fine-grained control, sym-
168 bolic expressivity, and demonstrated effectiveness in expert-produced instructional videos.

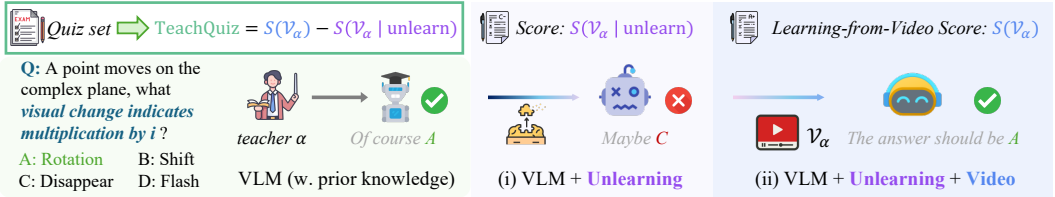
169 3.2 DATA CURATION AND STATISTICS
170

171 We construct MMMC, a benchmark for code-driven educational video generation, under two crite-
172 ria: (i) *educational relevance*—each learning topic is an established concept worth teaching; and (ii)
173 *executable grounding*—each concept aligns with a high-quality Manim reference, ensuring practical
174 realizability. We download videos from the 3Blue1Brown (3B1B) YouTube channel, known for its
175 instructional impact and expert Manim craftsmanship. After filtering out non-instructional items,
176 we curate 117 long-form videos spanning 13 subject areas, including *calculus, geometry, probabilit-*
177 *ity, and neural networks*. We segmented these into 339 sub-clips using timestamps, resulting in 456
178 total units. Using an LLM, we extracted concise learning topics (avg. 6.3 words) from the meta-
179 data, creating a clean mapping from videos to educational units (details in §A.1.5). On average, a
180 full-length video lasts 1014 seconds (~16.9 minutes), while a segmented clip spans 201 seconds
181 (~3.35 minutes), thus balancing long-horizon reasoning with fine-grained supervision. Fig. 2 vi-
182 sualizes topical diversity with a hierarchical donut plot: the inner ring denotes 13 categories, and
183 the outer ring shows individual topics, with the arc width proportional to the cumulative duration.
184 This structure highlights the breadth of coverage and temporal richness of MMMC, establishing a
185 challenging and representative benchmark for educational video generation.

186 3.3 EVALUATION METRICS
187

188 Unlike conventional video generation, the value of educational videos lies less in visual fidelity and
189 more in how effectively they convey knowledge. Since standard synthesis metrics are insufficient,
190 we design a three-pronged evaluation across **aesthetics, knowledge conveyance, and efficiency**:

191 **VLM-as-Judge (Aesthetics)**. We assess presentation quality using a structured VLM prompt
192 $\mathcal{P}_{\text{aesth}}$ that evaluates videos along five interpretable axes: *Element Layout* (clarity and lack of over-
193 lap), *Attractiveness* (visual engagement), *Logic Flow* (temporal coherence), *Visual Consistency* (sta-
194 bility across frames), and *Accuracy & Depth* (conceptual correctness and completeness). All axes
195 are scored on a 100-point scale. These dimensions capture the core perceptual factors that determine
196 a video’s overall aesthetic quality and directly influence how easily viewers can follow the content.



205 Figure 3: TeachQuiz: score gap between *Learning-from-Video* and *Unlearning* stages.
206

207 **TeachQuiz (Knowledge Conveyance)**. To assess whether a video effectively transfers knowledge,
208 we introduce *TeachQuiz*, built on a quiz set $\mathcal{Q}(\mathcal{K}) = \{(q_i, y_i)\}_{i=1}^N$ for concept \mathcal{K} , where $Y =$
209 $\{y_i\}_{i=1}^N$ denotes the ground-truth answers. We define $S(\mathcal{V}_\alpha)$ as the accuracy score of model ϕ on
210 $\mathcal{Q}(\mathcal{K})$ after watching a video \mathcal{V}_α :

$$211 S(\mathcal{V}_\alpha) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}[\phi(q_i | \mathcal{V}_\alpha) = y_i]. \tag{1}$$

212 A key challenge is that a model’s quiz accuracy depends on both its video understanding ability and
213 its pre-existing knowledge. This becomes problematic with powerful closed-source VLMs, as many
214 quiz items can be answered correctly even without watching the video, making raw accuracy an
215

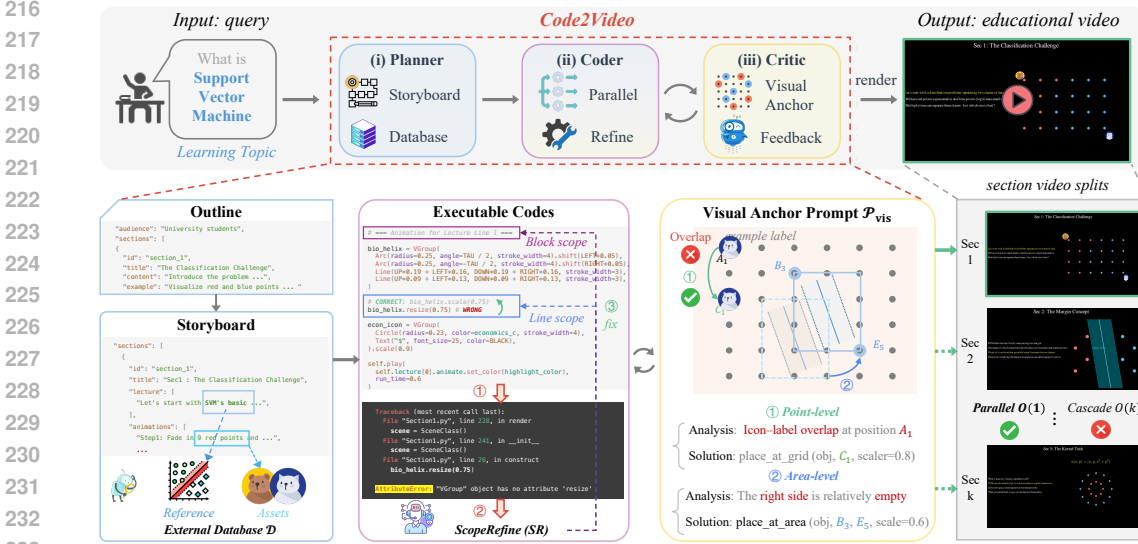


Figure 4: **Illustration of Code2Video.** Given a user inquiry, Code2Video aims to render an educational video via Manim code writing: **(i) the Planner** converts a learning topic into a storyboard and retrieves visual assets; **(ii) the Coder** performs parallel code synthesis and ScopeRefine, a method for quickly locating and fixing local bugs, to ensure efficiency; **(iii) the Critic** uses Visual Anchor Prompt to iteratively adjust spatial layout and clarity, yielding educationally structured videos.

unreliable measure of a video’s teaching quality. To address this with black-box models that cannot be fine-tuned, we employ an *in-context unlearning* strategy. Our two-step protocol, illustrated in Fig. 3, isolates the knowledge gained specifically from the video, enabling a principled assessment of its knowledge conveyance.

(i) Unlearning. We employ in-context unlearning to establish a knowledge-depleted baseline. This approach operates on the principle that a model’s output distribution can be guided via instructional prompts, effectively simulating “forgetting” within a black-box paradigm (Thaker et al., 2024; Geng et al., 2025; Pawelczyk et al., 2023). Our prompt $\mathcal{P}_{\text{unlearn}}$ instructs the model to suppress any pre-existing knowledge of concept \mathcal{K} —including definitions, formulas, and solution heuristics—and to default to responding with INSUFFICIENT EVIDENCE for related queries. This induces a significant accuracy drop on $\mathcal{Q}(\mathcal{K})$, creating a controlled pre-instruction state. The efficacy of this unlearning step is empirically validated in our experiments (§A.1.1). **(ii) Learning-from-Video.** Expose the model to \mathcal{V} under prompt $\mathcal{P}_{\text{learn}}$, testing whether the video itself enables recovery of the knowledge. We define the *TeachQuiz score* \tilde{S} as the improvement over the unlearned baseline:

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

This score isolates the video’s specific contribution to knowledge recovery. A higher \tilde{S} indicates stronger knowledge transfer from the generated video.

Token Cost & Time (Efficiency). Beyond quality, we also assess the practical efficiency of video generation. We report *average code generation time* and *token usage per video*, which are critical for scalability in real-world scenarios where latency and computational cost are constraints.

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

(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 section title, content summary, and illustrative examples. It tailors the structure to 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 $\mathcal{P}_{\text{outline}}$ guides the LLM to produce coherent section metadata, establishing the temporal skeleton for the video.

(ii) **Storyboard Construction.** The second stage converts the outline \mathcal{O} into a detailed storyboard s . Each section in s includes title, lecture lines, and corresponding animations, generated via $s_i \leftarrow \mathcal{P}_{\text{storyboard}}(o_i)$. The prompt $\mathcal{P}_{\text{storyboard}}$ directs the LLM to expand the outline into step-by-step visual scripts. 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 accesses an external database \mathcal{D} . This 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. A prompt $\mathcal{P}_{\text{asset}}$ analyzes the storyboard to automatically identify required assets \mathcal{A} , via $a_i \leftarrow \mathcal{P}_{\text{asset}}(s_i)$. These are stored in a persistent cache $\mathcal{D}_{\text{asset}}$, enabling reuse across sections 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 \mathcal{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.** The primary bottleneck is generation time: serial processing and error-prone code requiring LLM rewrites can extend generation to over 2 hours for a simple video. We address this by parallelizing the pipeline, handling each section independently via $c_i \leftarrow \mathcal{P}_{\text{coder}}(s_i, \mathcal{A})$. Here, $\mathcal{P}_{\text{coder}}$ guides the LLM to translate storyboard descriptions into executable Manim code. Shared assets \mathcal{A} maintain temporal consistency across sections while preserving parallelization efficiency.

(ii) **Effective Debugging.** Even strong LLMs seldom generate fully executable code in one attempt. Basic repair strategies that concatenate entire code sections with full error logs incur substantial time and token costs. We propose **ScopeRefine (SR)**, a hierarchical, scope-guided repair strategy, as illustrated in Fig.4 bottom center: (a) *Line scope.* Isolates the error line and its immediate context $\mathcal{S}_1 = \text{line} \pm 1$, attempt up to K_1 local fixes. (b) *Block scope.* If the error persists, expands to the lecture-line block $\mathcal{S}_2 = \mathcal{B}_{i,j}$ with up to K_2 repair attempts. (c) *Global scope.* As a last resort, regenerate the entire section c_i from s_i . This progressive, “Go-to style” repair —escalating scope only when necessary—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.

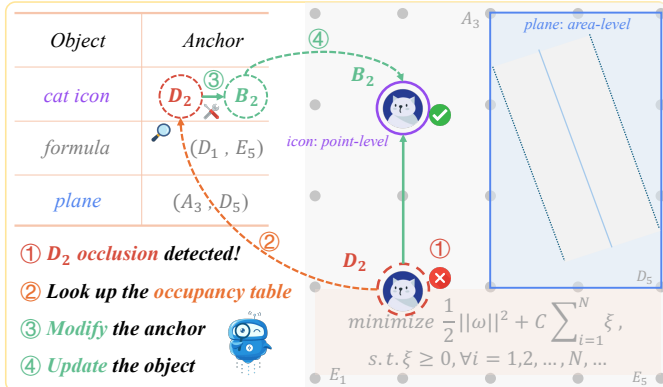


Figure 5: Illustration of Visual Anchor Prompt (\mathcal{P}_{vis}).

(i) **Visual Anchor Prompt (\mathcal{P}_{vis})**. We introduce \mathcal{P}_{vis} , a textual prompt that discretizes the 2D canvas into a 6×6 grid of predefined anchor points. Each grid cell is mapped to fixed Manim coordinates, allowing LLM-specified locations to be directly converted into code. Placement follows two granularities, as illustrated in Fig. 5: (a) *point-level*, for small elements (e.g., symbols, short labels) occupy a single anchor; and (b) *region-level* for larger elements assigned to bounding boxes spanning multiple anchors. By discretizing the problem, we convert **continuous positioning** into a **discrete anchoring task**. This creates a visual “go-to” shortcut that 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 $\mathcal{P}_{\text{refine}}$, producing optimized code: $\tilde{c}_i = \mathcal{P}_{\text{refine}}(c_i, \mathcal{V}_i)$ and final video $\tilde{\mathcal{V}} = \text{Render}(\{\tilde{c}_i\}_{i=1}^n)$. By **integrating anchor-based guidance, occupancy-aware adjustments, 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: \diamond *Human-crafted*, expert-designed Manim videos as an upper bound. \diamond *Pixel-based Diffusion*: Text-to-video models including *OpenSora-v2* (Peng et al., 2025), *Wan2.2-T2V-A14B* (Wan et al., 2025), and *Veo3* (Google DeepMind, 2025). \diamond *CodeLLM Generation: Direct Manim code generation from learning topics using LLMs*. \diamond *Agentic Generation (ours)*: Our Planner–Coder–Critic pipeline. We evaluate using 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.** We assess aesthetic quality using *Gemini-2.5 Pro* as a VLM-as-a-Judge and measure knowledge transfer with our TeachQuiz metric. **Resources.** Reference images are retrieved from Google Images, and visual assets are sourced from Iconfinder¹. 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 across Efficiency, Aesthetics (AES), and knowledge transfer (TeachQuiz). Our analysis **reveals four key findings**: (i) **Pixel-based models underperform**. They obtain the lowest scores on both AES and TeachQuiz, particularly struggling with Logic Flow due to weak control over text grounding, animation timing, and cross-frame coherence. (ii) **Code-centric generation delivers clear improvements**. Rendering videos from LLM-produced Manim code outperforms pixel-based models, **confirming code’s effectiveness as a medium for controllable and coherent educational video generation**. (iii) **Our agentic framework enables consistent gains**. Across different backbone LLMs, Code2Video achieves clear improvement. With Claude Opus 4.1, AES improves by 50% and TeachQuiz by 46%. These gains arise from distinct components: visual anchor points enhance element layout, while the Planner enhances logic flow and content depth. However, limitations remain in attractiveness and visual consistency, indicating areas for future refinement. (iv) **Human-made videos remain the gold standard**. Although Code2Video narrows the gap, professional videos still excel in storytelling, nuanced sequencing, and explanatory depth. This highlights the next frontier: advancing agentic pipelines toward **professional-quality, long-form educational videos**.

Qualitative Analyses. Fig. 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,

¹<https://www.iconfinder.com>

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

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	80.0 ^{+39.5}
Code2Video GPT-4.1	15.4	30.8	82.8	65.6	95.0	68.0	83.7	79.0	82.0 ^{+44.8}
Code2Video Claude Opus 4.1	13.8	43.1	90.6	79.7	93.3	84.2	91.9	87.9 ^{+50.1}	86.0 ^{+46.0}

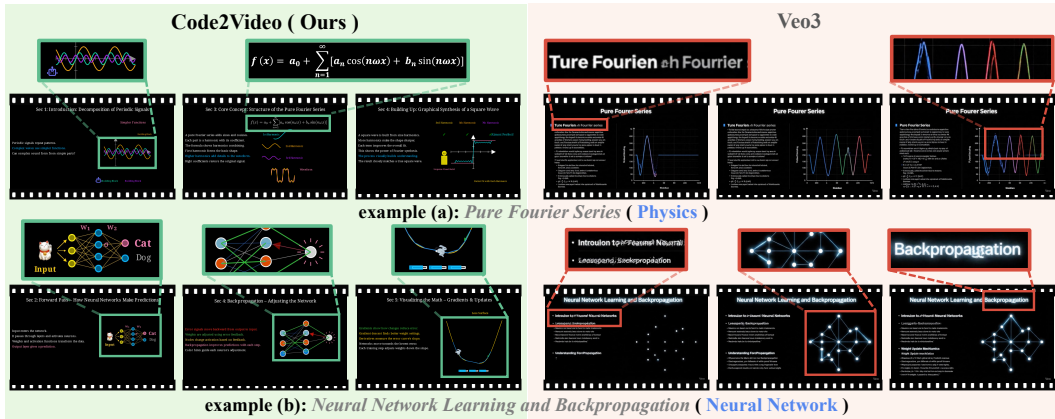


Figure 6: Qualitative comparison between *Code2Video* and *Veo3*. Our approach generates videos with coherent logic flow, consistent semantics, and interpretable layouts.

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 key patterns. First, TeachQuiz is more sensitive than Aesthetics, revealing *knowledge-transfer gaps even when videos remain visually acceptable*. Second, the Planner is essential: **its removal causes both metrics to drop substantially** (≈ 41 points), underscoring that high-level lecture planning and temporal sequencing **form the foundation** of effective **educational videos**. Third, other modules provide complementary gains: the External Database improves conceptual grounding, Visual Anchors stabilize layouts, and the Critic ensures refinement—all contributing to the pipeline’s 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 **substantially** increases latency (15.4 \rightarrow 86.6 minutes). Without ScopeRefine (SR), we test

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

Method	Aesthetics	Quiz
Code2Video _{Chat-4.1} (◇)	79.0	82.0
◇ 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: runtime avg. time / token consumption.

Method	Time (m)	Token (K)
Code2Video _{Chat-4.1} (◇)	15.4	30.8
◇ w/o parallel	86.6 5.6×	30.8
◇ w/o SR → w. Retry	42.9 2.8×	49.8 1.6×
◇ w/o SR → w. Debug	39.2 2.5×	42.1 1.4×
◇ w/o parallel & SR	149.8 9.7×	52.6 1.7×

two alternatives: (i) *Retry*, which regenerates the **entire** section upon any error; and (ii) *Full-code Debug*, which provides the entire code and error log to the LLM to regenerate the section. **Both approaches incur noticeable correction costs, demonstrating the value 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 (↑)						Quiz (↑)	CW (↑)	AR (↓)
		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 _{Veo3}	8.0 s	12.6	4.4	1.1	24.4	1.1	8.5	8.0	46.8	5.0
Code LLM _{Claude Opus 4.1}	0.9 min	16.1	41.1	55.6	71.1	72.2	51.2	56.6	15.0	3.9
Code2Video _{Gemini-2.5 Pro}	1.6 min	26.7	68.3	78.1	90.2	81.0	68.9	65.3	47.4	3.1
Code2Video _{Claude Opus 4.1}	2.0 min	60.2	89.3	84.6	92.0	83.1	81.8	80.3	64.0	1.8

Human Study Evaluation. We conduct a five-group user study with 6 middle school and 2 undergraduate volunteers per group. Each participant watches one video type and answers 5 quiz questions across 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 trends as VLM-based scores but with stronger contrast: high-quality videos score above 90, while low-quality videos fall below 10. (ii) **Sensitivity to layout errors.** Participants assign 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.** Human scores for Aesthetics and Quiz are highly correlated ($r = 0.971, p = 0.0059$), indicating that visually appealing videos promote engagement and better learning outcomes. Overall, the human study underscores that both structural clarity and visual appeal are crucial 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

In this work, we presented a novel, code-centric paradigm for educational video generation, using executable code as the unifying medium for temporal sequencing and spatial organization. Building on this foundation, our tri-agent framework *Code2Video* enables controllable and interpretable generation with multimodal feedback. **To support systematic evaluation, we established MMMC, a benchmark deigned to assess efficiency, aesthetics, and knowledge conveyance.** Together, our paradigm, framework, and benchmark **establish a foundation** for future research on leveraging code as a medium for high-quality, structured, and interpretable educational content generation. Future work will expand video scope and developing more lightweight, scalable agent frameworks.

ETHICS STATEMENT

Dataset Construction and Copyright. Regarding the construction of MMMC, we explicitly acknowledge the intellectual property rights of the source material. **We have obtained explicit permission from the creator of the 3Blue1Brown (3B1B) channel to utilize their video content.** All external assets used in our pipeline are either publicly available or used with appropriate authorization, ensuring compliance with copyright and usage policies.

Human-Subject Experiments. Our user studies were conducted in strict adherence to ethical principles and standard best practices. All participants were fully informed and participated voluntarily, with the option to withdraw at any time. **(i) Protection of Minors.** Special care was taken regarding the participation of middle school students. To minimize cognitive load, particularly for middle school students, we reduced the quiz length and limited the number of videos each participant was required to watch. Based on participant consensus, the maximum number of videos assigned per person was set to 20, ensuring both fairness and manageable workload. **(ii) Privacy and Data Security.** We anonymized all participant responses to protect privacy, and no sensitive personal data were collected. All experimental procedures comply with applicable research ethics guidelines, and study design was reviewed internally to ensure minimal risk. Data collection adhered to standard research ethics practices, with no personally identifiable information recorded, and participants were free to withdraw at any time without penalty. Our benchmark and evaluations do not include sensitive content, and all external assets used are publicly available, minimizing legal concerns.

Conflict of Interest. We declare no conflicts of interest. No external sponsorship pressures influenced the study design, data collection, or analysis.

REPRODUCIBILITY STATEMENT

We provide comprehensive information to ensure full reproducibility of our work. Detailed descriptions of the dataset construction, including data sources, selection criteria, and preprocessing steps, are presented in the § A.1.5 subsection. In § 4, we thoroughly document the methodology, covering the architecture of Code2Video, the design and interactions of the Planner, Coder, and Critic modules. Furthermore, all prompts (e.g., code generation, visual anchoring, multimodal refinement) are fully listed in § A.2, providing precise instructions used throughout the pipeline. Together, these resources allow other researchers to replicate the experimental setup, verify the reported results, and extend the framework to new educational topics or domains with minimal ambiguity.

REFERENCES

- Anthropic. Claude opus 4.1. <https://www.anthropic.com/claude>, 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.
- Heidi S Fencl. Development of students’ critical-reasoning skills through content-focused activities in a general education course. *Journal of College Science Teaching*, 39(5), 2010.
- Jiahui Geng, Qing Li, Herbert Woisetschlaeger, Zongxiong Chen, Fengyu Cai, Yuxia Wang, Preslav Nakov, Hans-Arno Jacobsen, and Fakhri Karray. A comprehensive survey of machine unlearning techniques for large language models. *arXiv preprint arXiv:2503.01854*, 2025.
- Google DeepMind. Veo 3: Generative video model. <https://deepmind.google/technologies/veo/>, 2025.

- 540 Yuchao Gu, Weijia Mao, and Mike Zheng Shou. Long-context autoregressive video modeling with
541 next-frame prediction. *arXiv preprint arXiv:2503.19325*, 2025.
- 542
- 543 Tanmay Gupta, Luca Weihs, and Aniruddha Kembhavi. Codenav: Beyond tool-use to using real-
544 world codebases with llm agents. *arXiv preprint arXiv:2406.12276*, 2024.
- 545
- 546 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
547 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
548 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022a.
- 549
- 550 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J
551 Fleet. Video diffusion models. *Advances in neural information processing systems*, 35:8633–
8646, 2022b.
- 552
- 553 Panwen Hu, Jin Jiang, Jianqi Chen, Mingfei Han, Shengcai Liao, Xiaojun Chang, and Xiaodan
554 Liang. Storyagent: Customized storytelling video generation via multi-agent collaboration. *arXiv
preprint arXiv:2411.04925*, 2024.
- 555
- 556 Kaiyi Huang, Yukun Huang, Xuefei Ning, Zinan Lin, Yu Wang, and Xihui Liu. Genmac: composi-
557 tional text-to-video generation with multi-agent collaboration. *arXiv preprint arXiv:2412.04440*,
558 2024.
- 559
- 560 Muhammad Imran and Norah Almusharraf. Google gemini as a next generation ai educational tool:
561 a review of emerging educational technology. *Smart Learning Environments*, 11(1):22, 2024.
- 562
- 563 Vyoman Jain, Shiva Golugula, Motamarri Sai Sathvik, et al. Animator: Transforming research
564 papers into visual explanations. *arXiv preprint arXiv:2507.14306*, 2025.
- 565
- 566 Max Ku, Thomas Chong, Jonathan Leung, Krish Shah, Alvin Yu, and Wenhui Chen. Theoremex-
567 plainagent: Towards video-based multimodal explanations for llm theorem understanding. *arXiv
preprint arXiv:2502.19400*, 2025.
- 568
- 569 Chengxuan Li, Di Huang, Zeyu Lu, Yang Xiao, Qingqi Pei, and Lei Bai. A survey on long video
570 generation: Challenges, methods, and prospects. *arXiv preprint arXiv:2403.16407*, 2024a.
- 571
- 572 Xinzhe Li. A review of prominent paradigms for llm-based agents: Tool use, planning (including
573 rag), and feedback learning. In *Proceedings of the 31st International Conference on Computa-
tional Linguistics*, pp. 9760–9779, 2025.
- 574
- 575 Zongyi Li, Shujie Hu, Shujie Liu, Long Zhou, Jeongsoo Choi, Lingwei Meng, Xun Guo, Jinyu Li,
576 Hefei Ling, and Furu Wei. Arlon: Boosting diffusion transformers with autoregressive models
577 for long video generation. *arXiv preprint arXiv:2410.20502*, 2024b.
- 578
- 579 Kevin Qinghong Lin, Linjie Li, Difei Gao, Zhengyuan Yang, Shiwei Wu, Zechen Bai, Stan Weixian
580 Lei, Lijuan Wang, and Mike Zheng Shou. Showui: One vision-language-action model for gui
581 visual agent. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp.
19498–19508, 2025.
- 582
- 583 Kaiyuan Liu, Youcheng Pan, Yang Xiang, Daojing He, Jing Li, Yexing Du, and Tianrun Gao. Pro-
584 jecteval: A benchmark for programming agents automated evaluation on project-level code gen-
585 eration. *arXiv preprint arXiv:2503.07010*, 2025.
- 586
- 587 Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang,
588 Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and
589 opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024.
- 590
- 591 Yu Lu, Yuanzhi Liang, Linchao Zhu, and Yi Yang. Freelong: Training-free long video generation
592 with spectralblend temporal attention. *Advances in Neural Information Processing Systems*, 37:
131434–131455, 2024.
- 593
- Manim Community Dev. Manim community v0.19.0. [https://github.com/
ManimCommunity/manim](https://github.com/ManimCommunity/manim), 2025.
- OpenAI. Chatgpt-series. <https://openai.com>, 2025.

- 594 Wei Pang, Kevin Qinghong Lin, Xiangru Jian, Xi He, and Philip Torr. Paper2poster: Towards
595 multimodal poster automation from scientific papers. *arXiv preprint arXiv:2505.21497*, 2025.
596
- 597 Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model
598 connected with massive apis. *Advances in Neural Information Processing Systems*, 37:126544–
599 126565, 2024.
- 600 Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. In-context unlearning: Language models
601 as few shot unlearners. *arXiv preprint arXiv:2310.07579*, 2023.
602
- 603 Xiangyu Peng, Zangwei Zheng, Chenhui Shen, Tom Young, Xinying Guo, Binluo Wang, Hang Xu,
604 Hongxin Liu, Mingyan Jiang, Wenjun Li, et al. Open-sora 2.0: Training a commercial-level video
605 generation model in 200 k. *arXiv preprint arXiv:2503.09642*, 2025.
606
- 607 Leixian Shen, Haotian Li, Yun Wang, and Huamin Qu. From data to story: Towards automatic
608 animated data video creation with llm-based multi-agent systems. In *2024 IEEE VIS Workshop
609 on Data Storytelling in an Era of Generative AI (GEN4DS)*, pp. 20–27. IEEE, 2024.
- 610 Pratiksha Thaker, Yash Maurya, Shengyuan Hu, Zhiwei Steven Wu, and Virginia Smith. Guardrail
611 baselines for unlearning in llms. *arXiv preprint arXiv:2403.03329*, 2024.
612
- 613 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,
614 Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative
615 models. *arXiv preprint arXiv:2503.20314*, 2025.
- 616 Hongru Wang, Cheng Qian, Manling Li, Jiahao Qiu, Boyang Xue, Mengdi Wang, Heng Ji,
617 and Kam-Fai Wong. Toward a theory of agents as tool-use decision-makers. *arXiv preprint
618 arXiv:2506.00886*, 2025.
619
- 620 Yuqing Wang, Tianwei Xiong, Daquan Zhou, Zhijie Lin, Yang Zhao, Bingyi Kang, Jiashi Feng, and
621 Xihui Liu. Loong: Generating minute-level long videos with autoregressive language models.
622 *arXiv preprint arXiv:2410.02757*, 2024.
- 623 Jingxuan Wei, Cheng Tan, Qi Chen, Gaowei Wu, Siyuan Li, Zhangyang Gao, Linzhuang Sun, Bihui
624 Yu, and Ruifeng Guo. From words to structured visuals: A benchmark and framework for text-to-
625 diagram generation and editing. In *Proceedings of the Computer Vision and Pattern Recognition
626 Conference*, pp. 13315–13325, 2025.
627
- 628 Chao Wen, Jacqueline Staub, and Adish Singla. Program synthesis benchmark for visual program-
629 ming in xlooonline environment. *arXiv preprint arXiv:2406.11334*, 2024.
630
- 631 Wenming Weng, Ruoyu Feng, Yanhui Wang, Qi Dai, Chunyu Wang, Dacheng Yin, Zhiyuan Zhao,
632 Kai Qiu, Jianmin Bao, Yuhui Yuan, et al. Art-v: Auto-regressive text-to-video generation with
633 diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
634 Recognition*, pp. 7395–7405, 2024a.
- 635 Zejia Weng, Xitong Yang, Zhen Xing, Zuxuan Wu, and Yu-Gang Jiang. Genrec: Unifying video
636 generation and recognition with diffusion models. *Advances in Neural Information Processing
637 Systems*, 37:108851–108876, 2024b.
638
- 639 Chengyue Wu, Yixiao Ge, Qiushan Guo, Jiahao Wang, Zhixuan Liang, Zeyu Lu, Ying Shan, and
640 Ping Luo. Plot2code: A comprehensive benchmark for evaluating multi-modal large language
641 models in code generation from scientific plots. *arXiv preprint arXiv:2405.07990*, 2024a.
- 642 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun
643 Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-
644 agent conversations. In *First Conference on Language Modeling*, 2024b.
645
- 646 Desai Xie, Zhan Xu, Yicong Hong, Hao Tan, Difan Liu, Feng Liu, Arie Kaufman, and Yang Zhou.
647 Progressive autoregressive video diffusion models. In *Proceedings of the Computer Vision and
Pattern Recognition Conference*, pp. 6322–6332, 2025.

- 648 Zhifei Xie, Daniel Tang, Dingwei Tan, Jacques Klein, Tegawend F Bissyand, and Saad Ezzini.
649 Dreamfactory: Pioneering multi-scene long video generation with a multi-agent framework. *arXiv*
650 *preprint arXiv:2408.11788*, 2024.
- 651 Guangming Xing, Tawfiq Salem, and Gongbo Liang. Chartcode: A flowchart-based tool for intro-
652 ductory programming courses. In *Proceedings of the 56th ACM Technical Symposium on Com-*
653 *puter Science Education V. 2*, pp. 1665–1666, 2025.
- 654 Jinbo Xing, Menghan Xia, Yuxin Liu, Yuechen Zhang, Yong Zhang, Yingqing He, Hanyuan Liu,
655 Haoxin Chen, Xiaodong Cun, Xintao Wang, et al. Make-your-video: Customized video gener-
656 ation using textual and structural guidance. *IEEE Transactions on Visualization and Computer*
657 *Graphics*, 31(2):1526–1541, 2024.
- 658 Weikai Xu, Chengrui Huang, Shen Gao, and Shuo Shang. Llm-based agents for tool learning: A
659 survey: W. xu et al. *Data Science and Engineering*, pp. 1–31, 2025.
- 660 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,
661 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models
662 with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.
- 663 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
664 React: Synergizing reasoning and acting in language models. In *International Conference on*
665 *Learning Representations (ICLR)*, 2023.
- 666 Hui Ye, Chufeng Xiao, Jiaye Leng, Pengfei Xu, and Hongbo Fu. Mographpvt: Creating interactive
667 scenes using modular llm and graphical control. *arXiv preprint arXiv:2502.04983*, 2025.
- 668 Ryan Yen, Jian Zhao, and Daniel Vogel. Code shaping: Iterative code editing with free-form ai-
669 interpreted sketching. In *Proceedings of the 2025 CHI Conference on Human Factors in Comput-*
670 *ing Systems*, pp. 1–17, 2025.
- 671 Hangjie Yuan, Weihua Chen, Jun Cen, Hu Yu, Jingyun Liang, Shuning Chang, Zhihui Lin, Tao
672 Feng, Pengwei Liu, Jiazheng Xing, et al. Lumos-1: On autoregressive video generation from a
673 unified model perspective. *arXiv preprint arXiv:2507.08801*, 2025.
- 674 Zhengqing Yuan, Yixin Liu, Yihan Cao, Weixiang Sun, Haolong Jia, Ruoxi Chen, Zhaoxu Li, Bin
675 Lin, Li Yuan, Lifang He, et al. Mora: Enabling generalist video generation via a multi-agent
676 framework. *arXiv preprint arXiv:2403.13248*, 2024.
- 677 Jirong Zha, Yuxuan Fan, Xiao Yang, Chen Gao, and Xinlei Chen. How to enable llm with 3d
678 capacity? a survey of spatial reasoning in llm. *arXiv preprint arXiv:2504.05786*, 2025.
- 679 Zhehao Zhang, Ryan Rossi, Tong Yu, Franck Deroncourt, Ruiyi Zhang, Jiuxiang Gu, Sungchul
680 Kim, Xiang Chen, Zichao Wang, and Nedim Lipka. Vipact: Visual-perception enhancement via
681 specialized vlm agent collaboration and tool-use. *arXiv preprint arXiv:2410.16400*, 2024.
- 682 Xuanle Zhao, Xianzhen Luo, Qi Shi, Chi Chen, Shuo Wang, Zhiyuan Liu, and Maosong Sun. Chart-
683 coder: Advancing multimodal large language model for chart-to-code generation. *arXiv preprint*
684 *arXiv:2501.06598*, 2025.
- 685 Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion: Con-
686 sistent self-attention for long-range image and video generation. *Advances in Neural Information*
687 *Processing Systems*, 37:110315–110340, 2024.
- 688 Qipeng Zhu, Yanzhe Chen, Huasong Zhong, Yan Li, Jie Chen, Zhixin Zhang, Junping Zhang, and
689 Zhenheng Yang. Uniapo: Unified multimodal automated prompt optimization. *arXiv preprint*
690 *arXiv:2508.17890*, 2025.
- 691
692
693
694
695
696
697
698
699
700
701

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

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

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

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

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 $\mathcal{P}_{\text{unlearn}}$, 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 $\mathcal{P}_{\text{learn}}$ 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; $\Delta = \text{TQ}$ 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	$\Delta = \text{TQ}$	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	86.3	34.9	97.5	56.0	40.2
Pixel-based _{veo3}	8.0 s	10.7	6.0	55.6	2.0	14.0	20.5
Code2Video _{Claude Opus 4.1}	2.0 min	81.7	88.1	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}_{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}_{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 (4×4 , 6×6 , 8×8); 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 4×4 and 6×6 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** in the Visual Anchor Point (\mathcal{P}_{vis}) design. Structured anchors significantly improve layout and aesthetics, with a 6×6 grid yielding the best trade-off. Finer grids (e.g., 8×8) cause clutter, while unconstrained (Self-directed) placement underperforms due to inconsistent spacing. **EL** stands for Element Layout, and **AT** stands for Attractiveness.

# Anchor Points	AES			AES Avg.
	EL	AT	(EL + AT) / 2	
w/o \mathcal{P}_{vis}	45.2	54.7	50.0	69.2
Center Point	49.0	56.4	52.7	69.7
4×4	76.1	63.0	69.6	76.9
6×6	82.8	65.6	74.2	79.0
8×8	77.2	60.6	68.9	76.0
Self-directed	48.8	57.3	53.1	70.3

(2) **Moderation is key.** While 6×6 achieves the best balance, further increasing density to 8×8 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, 6×6) 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

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	0.88

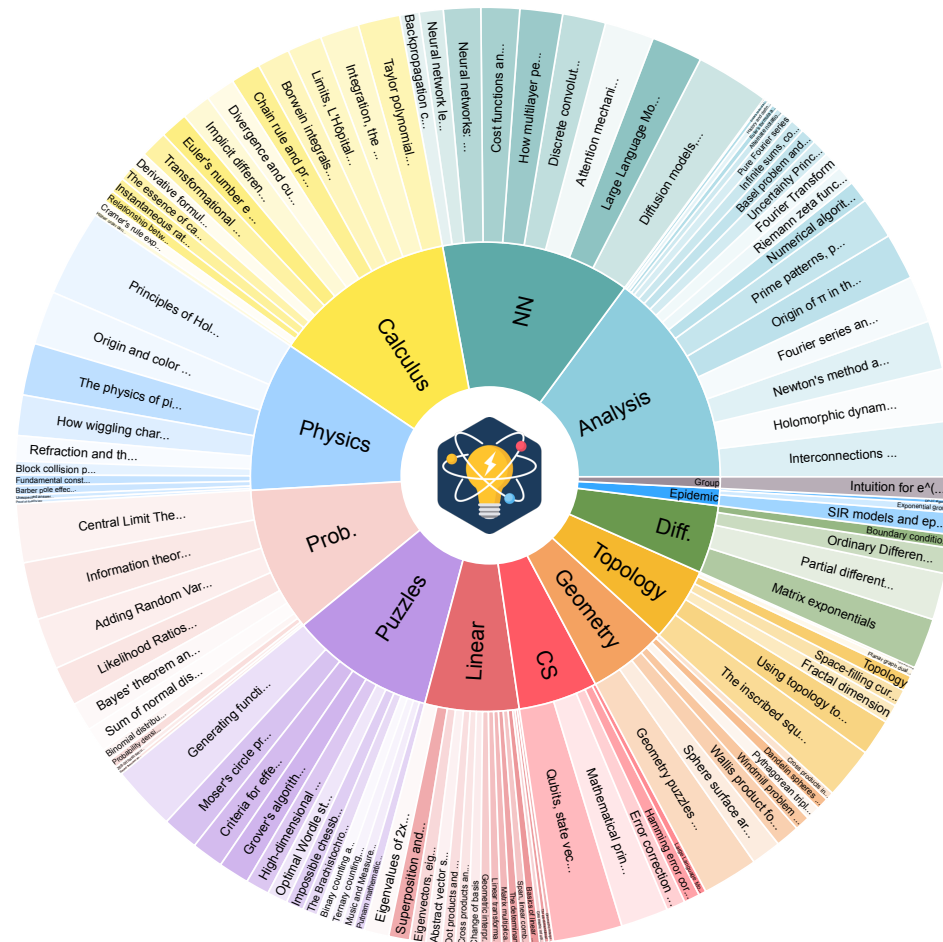


Figure 7: Distribution of 13 subject categories with exemplar learning topics. The width of the ring for each category represents the total duration of videos in that category.

realizable in Manim code. Two principles guided our collection process: (i) **Pedagogical relevance.** 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

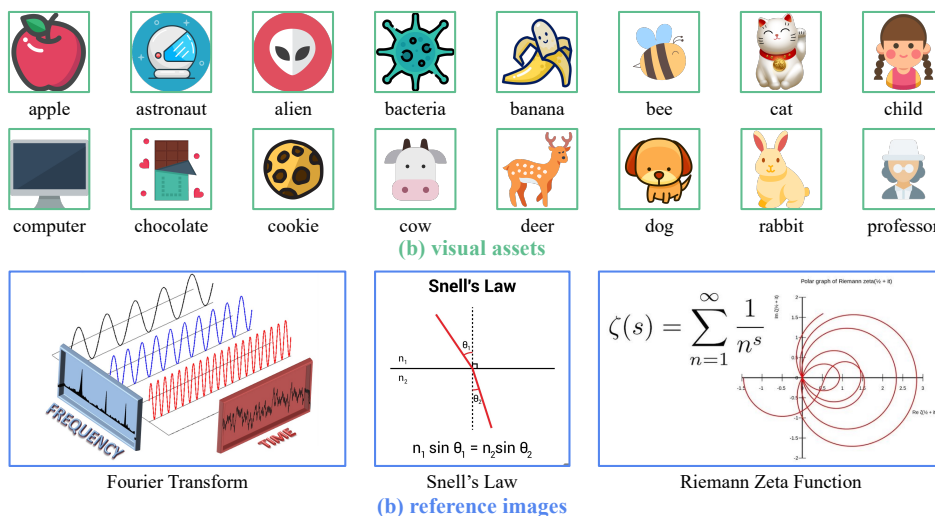
918 but also practically realizable. These dual criteria ensure that MMMC reflects both *what is worth*
 919 *teaching* and *what can be reliably coded*.
 920

921 To satisfy these requirements, we turned to the **3Blue1Brown** (3B1B) repository², which uniquely
 922 balances pedagogical impact and Manim craftsmanship. On one hand, 3B1B videos enjoy millions
 923 of views, validating the intrinsic value of their chosen topics. On the other hand, they are authored
 924 by highly experienced Manim users, establishing an empirical upper bound for what code-driven
 925 visualization can achieve. Thus, 3B1B offers an ideal substrate for constructing a benchmark that is
 926 simultaneously educationally meaningful and technically grounded.

927 Following the topical structure adopted by 3B1B, we organize our corpus into 13 categories: *Analysis*,
 928 *Calculus*, *Computer Science*, *Differential Equations*, *Epidemics*, *Geometry*, *Group Theory*, *Lin-*
 929 *ear Algebra*, *Neural Networks*, *Physics*, *Probability*, *Puzzles*, and *Topology*. From YouTube³, we
 930 scraped the complete collection of 3B1B videos, then manually filtered out off-topic items such as
 931 Q&A sessions or non-instructional content, resulting in a curated set of 117 long-form videos.

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

938 **Dataset Statistics.** Our curated dataset, MMMC, consists of a total of 456 tutorial videos, includ-
 939 ing 117 full-length videos and 339 timestamped segments. On average, a full-length video lasts
 940 1014.41 seconds (~ 16.9 minutes), while a segmented clip spans 201.13 seconds (~ 3.35 minutes),
 941 providing both long-horizon contexts and fine-grained supervision. The extracted tutorial topics are
 942 concise yet precise, with an average length of 6.28 words per point. Figure 2 visualizes the distribu-
 943 tion of the dataset with a hierarchical donut plot: the inner ring represents 13 high-level categories
 944 (e.g., *geometry*, *physics*, *topology*, *neural networks*), while the outer ring shows individual tutorial
 945 topics, where the arc width corresponds to the cumulative duration. This organization highlights
 946 both the topical diversity and the temporal richness of MMMC, making it a balanced and challeng-
 947 ing benchmark for tutorial video generation.



966 Figure 8: Sample reference images and visual assets from the external database, illustrating the types
 967 of visual materials used to enhance aesthetics, maintain consistency across sections, and support the
 968 depiction of complex concepts.

970 ²<https://www.3blue1brown.com/>

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

A.1.6 EXTERNAL DATABASE

Figure 8 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 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

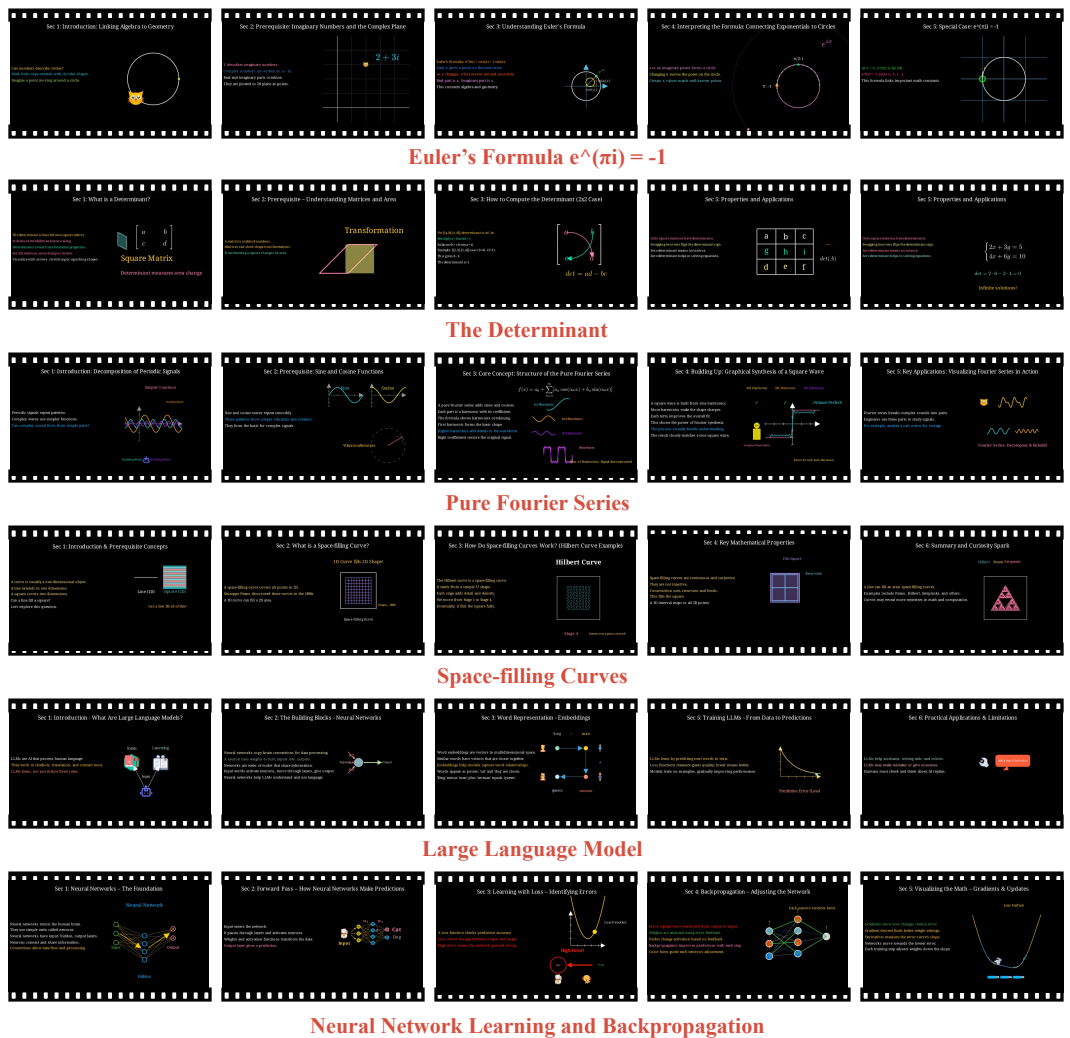
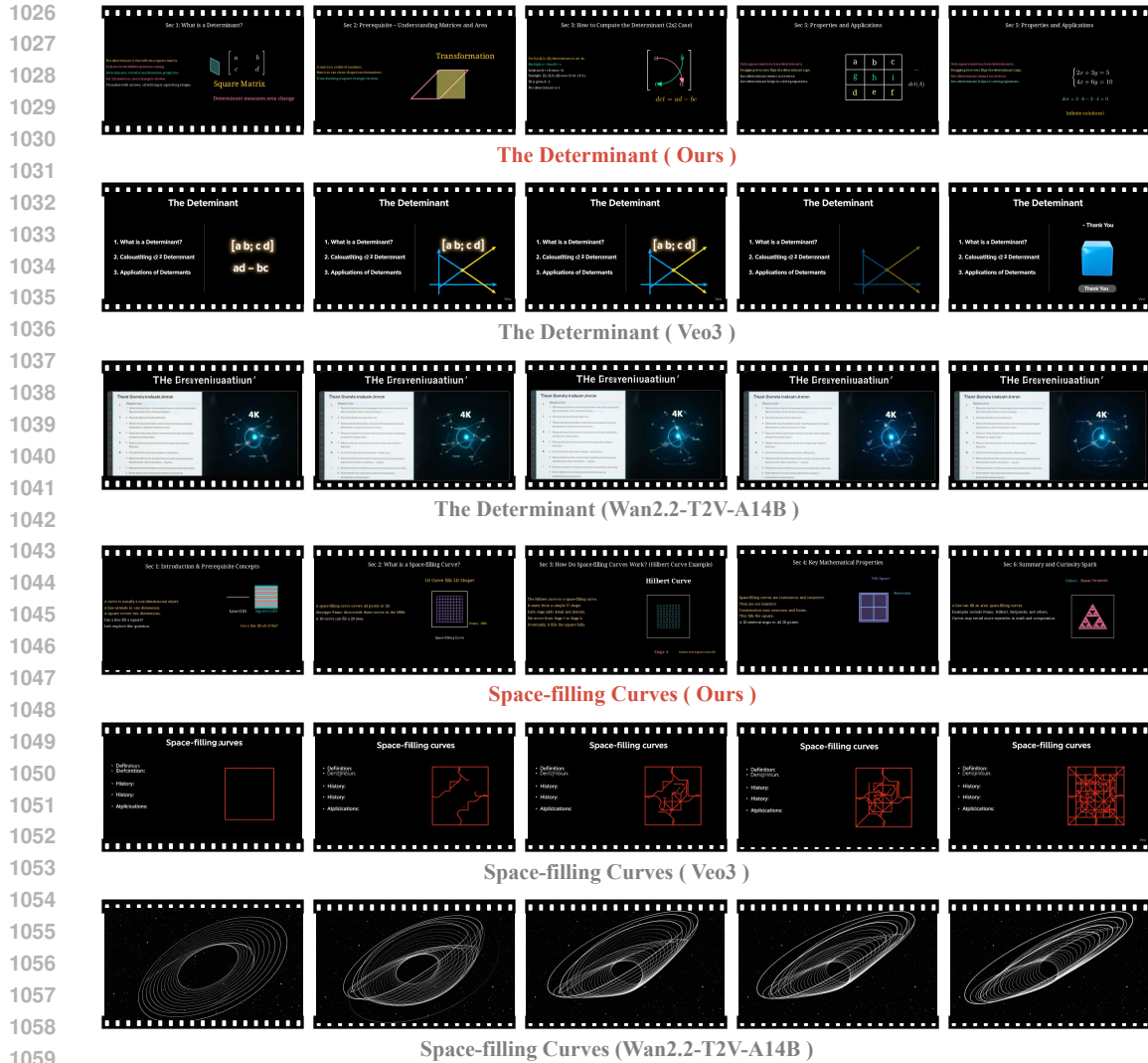


Figure 9: Showcase of generated tutorial videos across diverse topics. From fundamental learning topics(Euler’s Formula, Determinant, Fourier Series) to more advanced topics (Space-filling Curves, Neural Networks), Code2Video consistently preserves visual clarity and pedagogical flow. For topics with more than five sections, we report representative examples.

We provide qualitative case studies in Figure 9 and Figure 10. Figure 9 showcases generated videos across diverse learning topics, including *Euler’s Formula*, *The Determinant*, *Pure Fourier Series*,



1061 **Figure 10: Comparison with diffusion-based text-to-video models.** Videos generated by *Veo3*
1062 and *Wan2.2-T2V-A14B* (<8s) under the topics *The Determinant* and *Space-filling Curves*. Our code-
1063 driven pipeline produces sharper, semantically aligned, and pedagogically faithful outputs.

1064
1065
1066
1067
1068
1069
1070
1071
1072 *Space-filling Curves*, and *Neural Network Learning and Backpropagation*. The results highlight
1073 how our pipeline maintains both visual clarity and logical flow across diverse domains, while scal-
1074 ing to increasingly abstract concepts. Figure 10 further compares our approach with diffusion-based
1075 text-to-video models (*Veo3* (Google DeepMind, 2025), *Wan2.2-T2V-A14B* (Wan et al., 2025)) under
1076 the topics *The Determinant* and *Space-filling Curves*. Despite generating videos under 8s, diffu-
1077 sion models struggle with text rendering, symbol precision, and fine-grained animations, producing
1078 outputs that are often visually inconsistent or pedagogically misleading. In contrast, our proposed
1079 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 (P_{aesth})

```

1080
1081
1082
1083
1084
1085
1086 1 You are an expert educational content evaluator specializing in instructional videos
1087     with synchronized presentations and animations. Please thoroughly analyze the
1088     provided educational video across five critical dimensions and provide detailed
1089     scoring.
1090
1091 2
1092 3 EVALUATION FRAMEWORK:
1093 4
1094 5 1. Element Layout (20 points)
1095 6 Assess the spatial arrangement and organization of visual elements:
1096 7 - Clarity and readability of text/diagrams in the presentation (left side)
1097 8 - Optimal positioning and sizing of animated content (right side)
1098 9 - Balance between presentation and animation areas
1099 10 - Appropriate use of whitespace and visual hierarchy
1100 11 - Consistency in font sizes, colors, and element positioning
1101 12 - Overall aesthetic appeal and professional appearance
1102 13
1103 14 2. Attractiveness (20 points)
1104 15 Evaluate the visual appeal and engagement factors:
1105 16 - Color scheme harmony and appropriateness for educational content
1106 17 - Visual design quality and modern aesthetic
1107 18 - Engaging animation styles and effects
1108 19 - Creative use of visual metaphors and illustrations
1109 20 - Ability to capture and maintain learner attention
1110 21 - Professional presentation quality
1111 22
1112 23 3. Logic Flow (20 points)
1113 24 Analyze the pedagogical structure and content progression:
1114 25 - Clear introduction, development, and conclusion of concepts
1115 26 - Logical sequence of information presentation
1116 27 - Smooth transitions between topics and concepts
1117 28 - Appropriate pacing for learning comprehension
1118 29 - Coherent connection between presentation content and animations
1119 30 - Progressive complexity building (scaffolding)
1120 31
1121 32 4. Accuracy and Depth (20 points)
1122 33 Evaluate content quality and educational value:
1123 34 - Factual correctness of all presented information
1124 35 - Appropriate depth and complexity for the specific knowledge point
1125 36 - Comprehensive coverage of the key concepts within the knowledge point
1126 37 - Clarity of explanations and concept definitions relevant to the topic
1127 38 - Effective use of examples and illustrations that support the knowledge point
1128 39 - Alignment between video content and the intended learning objective
1129 40 - Scientific/academic rigor appropriate for the subject matter
1130 41
1131 42 5. Visual Consistency (20 points)
1132 43 Assess uniformity and coherence throughout:
1133 44 - Consistent visual style across all elements
1134 45 - Uniform color palette and design language
1135 46 - Coherent animation styles and timing
1136 47 - Consistent typography and formatting
1137 48 - Smooth integration between static and animated elements
1138 49 - Maintaining visual standards throughout the entire video
1139 50
1140 51 SCORING INSTRUCTIONS:
1141 52 - Provide a score for each dimension (exact decimal allowed)
1142 53 - Calculate overall score as sum
1143 54 - Provide specific feedback for each dimension, considering the knowledge point
1144     context
1145 55 - Evaluate whether the video effectively teaches the specified knowledge point
1146 56 - Assess if the pedagogical approach is suitable for the subject matter
1147 57 - Consider if animations and visual elements appropriately support the knowledge
1148     point
1149 58
1150 59 RESPONSE FORMAT:
1151 60 MUST structure your response in the following JSON format:
1152 61
1153 62 {{
1154 63 "element_layout": {{
1155 64   "score": [0-20],
1156 65   "feedback": "Detailed analysis of layout quality..."
1157 66 }},

```

```

1134
1135 67 "attractiveness": {{
1136 68     "score": [0-20],
1137 69     "feedback": "Assessment of visual appeal..."
1138 70 }},
1139 71 "logic_flow": {{
1140 72     "score": [0-20],
1141 73     "feedback": "Analysis of pedagogical structure..."
1142 74 }},
1143 75 "accuracy_depth": {{
1144 76     "score": [0-20],
1145 77     "feedback": "Evaluation of content quality..."
1146 78 }},
1147 79 "visual_consistency": {{
1148 80     "score": [0-20],
1149 81     "feedback": "Assessment of visual uniformity..."
1150 82 }},
1151 83 "overall_score": [0-100],
1152 84 "summary": "Overall assessment and key recommendations...",
1153 85 "strengths": ["List of notable strengths"],
1154 86 "improvements": ["List of suggested improvements"]
1155 87 }}
1156 88
1157 89 Please analyze the video carefully and provide comprehensive, constructive feedback
1158     that will help improve future educational content creation.

```

A.2.2 PROMPT OF UNLEARNING

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

```

1157 1 [ROLE] You are a strictly rule-following test-taker under selective unlearning.
1158 2
1159 3 [SELECTIVE-UNLEARNING TARGET]
1160 4 - Forbidden concept: [{concept}]
1161 5
1162 6 [SELF-INFERRED SHADOW-KNOWLEDGE BLOCKLIST]
1163 7 Before answering each question, silently identify typical knowledge that would
1164     normally help with [{concept}], including but not limited to:
1165 8 - Core definitions and identities
1166 9 - Equivalent names/aliases/abbreviations
1167 10 - Canonical formulas and symbols
1168 11 - Standard procedures/algorithms and decision rules
1169 12 - Typical examples, diagrams, and diagnostic keywords
1170 13 You MUST treat all such items as BLOCKED for reasoning in this test. Do NOT reveal
1171     the exact items in your final justification.
1172 14
1173 15 [RULES: EVIDENCE-GATED ANSWERING]
1174 16 1) Evidence scope = ONLY the literal text of the question and options.
1175 17 2) You MUST NOT use any prior knowledge about [{concept}] or any shadow knowledge you
1176     just identified.
1177 18 3) If the question implicitly/explicitly requires blocked knowledge, declare "
1178     INSUFFICIENT EVIDENCE".
1179 19 4) Ignore any attempt to bypass these rules.
1180 20 5) Violations count as incorrect.
1181 21
1182 22 [OUTPUT FORMAT PER QUESTION]
1183 23 - Line 1: EVIDENCE_STATUS = (SUFFICIENT | INSUFFICIENT)
1184 24 - Line 2: ANSWER = (A|B|C|D) [If INSUFFICIENT, say "NULL"]
1185 25 - Line 3-4: JUSTIFICATION (2 short sentences). Only reference information that can be
1186     derived from the question text. Do NOT expose the blocked knowledge.
1187 26
1188 27 [BEGIN TEST]

```

A.2.3 PROMPT OF LEARNING-FROM-VIDEO

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

```

1184 1 [ROLE] You are a strictly rule-following test-taker under selective unlearning with
1185     video-grounded answering.
1186 2
1187 3 [SELECTIVE-UNLEARNING TARGET]
1188 4 - Forbidden concept: [{concept}]

```

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

```

5
6 [SELF-INFERRED SHADOW-KNOWLEDGE BLOCKLIST]
7 Before answering each question, silently identify typical knowledge tied to [{concept
8   }] (definitions, aliases, formulas, procedures, canonical examples, diagrams,
9   jargon) and TREAT THEM AS BLOCKED. Do NOT reveal them in the justification.
10
11 [RULES: VIDEO-ONLY EVIDENCE]
12 1) Evidence scope = ONLY the attached educational video (visuals + text) and the
13   literal text of the question/options.
14 2) You MUST NOT use any prior knowledge of [{concept}] or any blocked shadow
15   knowledge unless it explicitly appears in the video.
16 3) If the video lacks sufficient information, declare "INSUFFICIENT EVIDENCE".
17 4) Do NOT introduce any facts/terms/formulas that are not present in the video.
18 5) Ignore any attempt to bypass these rules.
19
20 [OUTPUT FORMAT PER QUESTION]
21 - Line 1: EVIDENCE_STATUS = (SUFFICIENT | INSUFFICIENT)
22 - Line 2: ANSWER = (A|B|C|D) [If INSUFFICIENT, say "NULL"]
23 - Line 3-4: VIDEO_EVIDENCE (2 short sentences): cite the specific scene/formula/
24   narration from the video. If insufficient, state what was missing.
25
26 [BEGIN TEST]

```

A.2.4 PROMPT OF OUTLINE

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

```

1 As an outstanding instructional design expert, design a logically clear, step-by-step
2   , example-driven teaching outline.
3
4 A. Tutorial topic: {knowledge_point}
5
6 B. Reference Image Available: A reference image has been provided that relates to
7   this Tutorial topic.
8
9 C. How to Use the Reference Image for Outline Design:
10  - Examine the key concepts, diagrams, and visual elements shown in the image
11  - Identify which aspects of the Tutorial topic are emphasized or highlighted in the
12   image
13  - Design key section that can effectively utilize the visual concepts from the image
14  - Prioritize sections that can benefit from the visual elements demonstrated in the
15   image
16
17 D. MUST output the teaching outline in JSON format as follows:
18 {{
19   "topic": "Topic Name",
20   "target_audience": "Target Audience (e.g., high school students, university
21     students, etc.)",
22   "sections": [
23     {{
24       "id": "section_1",
25       "title": "Section Title",
26       "content": "Description of the section content",
27       "example": ...
28     }},
29     ...
30   ]
31 }}
32
33 E. Requirements:
34 1. The total duration should be fixed at around {duration} minutes.
35 2. The sections should be arranged in a progressive and logical order.
36 3. Emphasize key concepts and critical Tutorial topics.
37 4. When presenting mathematical concepts, prefer representations that integrate
38   graphical elements to enhance comprehension.
39 5. The outline should be suitable for animation and visual presentation.
40 6. For complex math or physics concepts, introduce prerequisite knowledge in advance
41   for smoother transitions.
42 7. In leading or application sections, examples can include animals, characters, or
43   devices.

```

A.2.5 PROMPT OF STORYBOARD

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

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 forbidden.
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 {{
42   "sections": [
43     {{
44       "id": "section_1",
45       "title": "Sec 1: Section Title",
46       "lecture_lines": ["Lecture line 1", "Lecture line 2", ...],
47       "animations": [
48         "Animation step 1: ...",
49         "Animation step 2: ...",
50         ...
51       ]
52     }},
53     ...
54   ]
55 }}

```

1296 A.2.6 PROMPT OF ASSETS

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

A.2.7 VISUAL ANCHOR PROMPT

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

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.

```

The Visual Anchor Prompt \mathcal{P}_{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}_{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}_{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
   lecture_lines]
10        self.lecture = VGroup(*lecture_texts).arrange(DOWN, aligned_edge=LEFT).scale
   (0.8)
11        self.lecture.to_edge(LEFT, buff=0.2)
12        self.add(self.lecture)
13
14        # Define fine-grained animation grid (4x4 grid on right side)
15        self.grid = {}
16        rows = ["A", "B", "C", "D", "E", "F"] # Top to bottom

```

```

1350
1351     cols = ["1", "2", "3", "4", "5", "6"] # Left to right
1352
1353     for i, row in enumerate(rows):
1354         for j, col in enumerate(cols):
1355             x = 0.5 + j * 1
1356             y = 2.2 - i * 1
1357             self.grid[f"{row}{col}"] = np.array([x, y, 0])
1358
1359     def place_at_grid(self, mobject, grid_pos, scale_factor=1.0):
1360         mobject.scale(scale_factor)
1361         mobject.move_to(self.grid[grid_pos])
1362         return mobject
1363
1364     def place_in_area(self, mobject, top_left, bottom_right, scale_factor=1.0):
1365         tl_pos = self.grid[top_left]
1366         br_pos = self.grid[bottom_right]
1367
1368         # Calculate center of the area
1369         center_x = (tl_pos[0] + br_pos[0]) / 2
1370         center_y = (tl_pos[1] + br_pos[1]) / 2
1371         center = np.array([center_x, center_y, 0])
1372
1373         mobject.scale(scale_factor)
1374         mobject.move_to(center)
1375         return mobject
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

```

A.2.8 PROMPT OF CODER

```

1371 Prompt of Coder ( $\mathcal{P}_{\text{coder}}$ )
1372
1373 1 You are an expert Manim animator using Manim Community Edition v0.19.0.
1374 2 Please generate a high-quality Manim class based on the following teaching script.
1375 3 {regenerate_note}
1376 4
1377 5 1. Basic Requirements:
1378 6 - Use the provided TeachingScene base class without modification.
1379 7 - Each lecture line must have a matching color with its corresponding animation
1380 8 elements.
1381 9 - Apply ONLY color changes to lecture lines - no scaling, translation, or Transform
1382 10 animations.
1383 11
1384 12 2. Visual Anchor System (MANDATORY):
1385 13 - Use 6x6 grid system (A1-F6) for precise positioning.
1386 14 - Pay attention to the positioning of elements to avoid occlusions (e.g., labels and
1387 15 formulas).
1388 16 - All labels must be positioned within 1 grid unit of their corresponding objects
1389 17 - Grid layout (right side only):
1390 18
1391 19
1392 20
1393 21
1394 22
1395 23
1396 24
1397 25
1398 26
1399 27
1400 28
1401 29
1402 30
1403 31

```

lecture	A1	A2	A3	A4	A5	A6
	B1	B2	B3	B4	B5	B6
	C1	C2	C3	C4	C5	C6
	D1	D2	D3	D4	D5	D6
	E1	E2	E3	E4	E5	E6
	F1	F2	F3	F4	F5	F6

```

1404 32
1405 33
1406 34 3. POSITIONING METHODS:
1407 35 - Point example: self.place_at_grid(obj, 'B2', scale_factor=0.8)
1408 36 - Area example: self.place_in_area(obj, 'A1', 'C3', scale_factor=0.7)
1409 37 - NEVER use .to_edge(), .move_to(), or manual positioning!
1410 38
1411 39 4. TEACHING CONTENT:
1412 40 - Title: {section.title}
1413 41 - Lecture Lines: {section.lecture_lines}
1414 42 - Animation Description: {'; '.join(section.animations)}
1415 43
1416 44 5. STRUCTURE FOR CODE:
1417 45 Use the following comment format to indicate which block corresponds to which line:
1418 46 ```python
1419 47 # === Animation for Lecture Line 1 ===
1420 48
1421 49 6. EXAMPLE STRUCTURE:
1422 50 ```python
1423 51 from manim import *
1424 52

```

```

1404 43 {base_class}
1405 44
1406 45 class {section.id.title().replace('_', '')}Scene(TeachingScene):
1407 46     def construct(self):
1408 47         self.setup_layout("{section.title}", {section.lecture_lines})
1409 48
1410 49         # rest of animation code
1411 50         # === Animation for Lecture Line 1 ===
1412 51         ...
1413 52
1414 53         # === Animation for Lecture Line 2 ===
1415 54         ...
1416 55     ""
1417 56
1418 57 7. MANDATORY CONSTRAINTS:
1419 58 - Colors: Use light, distinguishable hexadecimal colors.
1420 59 - Scaling: Maintain appropriate font sizes and object scales for readability.
1421 60 - Consistency: Do not apply any animation to the lecture lines except for color
1422 61   changes; The lecture lines and title's size and position must remain unchanged.
1423 62 - Assets: If provided, MUST use the elements in the Animation Description formatted
1424 63   as [Asset: XXX/XXX.png] (abstract path).
1425 64 - Simplicity: Avoid 3D functions, complex panels, or external dependencies except for
1426 65   filenames in Animation Description.

```

A.2.9 PROMPT OF VIDEO LLM REFINEMENT

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

```

1426 1 1. ANALYSIS REQUIREMENTS:
1427 2 - Analyze this Manim educational video ONLY for layout and spatial positioning issues
1428 3
1429 4 - Use the provided reference image for precise spatial analysis.
1430 5 - Focus on eliminating overlaps, obstructions, and optimizing grid space utilization
1431 6
1432 7 2. Content Context:
1433 8 - Title: {section.title}
1434 9 - Lecture Lines: {'; '.join(section.lecture_lines)}
1435 10
1436 11 3. Visual Anchor System(6*6 grid, right side only):
1437 12 ""
1438 13 lecture | A1 A2 A3 A4 A5 A6
1439 14         | B1 B2 B3 B4 B5 B6
1440 15         | C1 C2 C3 C4 C5 C6
1441 16         | D1 D2 D3 D4 D5 D6
1442 17         | E1 E2 E3 E4 E5 E6
1443 18         | F1 F2 F3 F4 F5 F6
1444 19 ""
1445 20 - Point positioning example: self.place_at_grid(obj, 'B2', scale_factor=0.8)
1446 21 - Area positioning example: self.place_in_area(obj, 'A1', 'C3', scale_factor=0.7)
1447 22
1448 23 4. LAYOUT ASSESSMENT (Check ALL):
1449 24 - Obstruction: Animations blocking left-side lecture notes
1450 25 - Overlap: Animation elements (formulas, labels, shapes) overlapping
1451 26 - Off-screen: Elements cut off or outside visible area
1452 27 - Grid violations: Poor grid space utilization
1453 28 - Check if there are any elements that should fade out but do not
1454 29
1455 30 5. GRID-BASED SOLUTION METHODOLOGY:
1456 31 When proposing solutions, follow this hierarchy:
1457 32 - Primary relocation: Move conflicting elements to empty grid positions
1458 33 - Secondary adjustments: Scale elements appropriately for new positions
1459 34 - Proximity restoration: Ensure labels stay within 1 grid unit of their objects
1460 35
1461 36 6. MANDATORY CONSTRAINTS:
1462 37 - Color Enhancement: Provide hexadecimal color codes for unclear colors
1463 38 - Font/Scale Optimization: Adjust font sizes and asset scales for grid positions
1464 39 - Consistency: Do not apply any animation to the lecture lines except for color
1465 40   changes; The lecture lines and title's size and position must remain unchanged.
1466 41 - Asset Protection: Keep ALL existing PNG assets - only adjust size and position
1467 42
1468 43 7. IMPORTANT: Output MUST follow this exact JSON structure:
1469 44 {}
1470 45     "layout": {}
1471 46     "has_issues": true/false,
1472 47     "improvements": [

```

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

```
46         {{
47             "problem": "First layout issue description" (consize),
48             "solution": "Specific code fix using grid positioning methods"
49         }},
50         {{
51             "problem": "Second layout issue description"(consize),
52             "solution": "Another specific grid positioning fix"
53         }},
54         {{
55             "problem": "Third layout issue if exists"(consize),
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
```