

# Speech-Synchronized Whiteboard Generation via VLM-Driven Structured Drawing Representations

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Creating whiteboard-style educational videos demands precise*  
002 *coordination between freehand illustrations and spoken*  
003 *narration, yet no existing method addresses this multimodal*  
004 *synchronization problem with structured, reproducible*  
005 *drawing representations. We present the first dataset of*  
006 *24 paired Excalidraw demonstrations with narrated audio,*  
007 *where every drawing element carries millisecond-precision*  
008 *creation timestamps spanning 8 STEM domains. Using this*  
009 *data, we study whether a vision-language model (Qwen2-*  
010 *VL-7B), fine-tuned via LoRA, can predict full stroke se-*  
011 *quences synchronized to speech from only 24 demonstra-*  
012 *tions. Our topic-stratified five-fold evaluation reveals*  
013 *that timestamp conditioning significantly improves tempo-*  
014 *ral alignment over ablated baselines, while the model ge-*  
015 *neralizes across unseen STEM topics. We discuss trans-*  
016 *ferability to real classroom settings and will release our*  
017 *dataset and code upon acceptance to support future re-*  
018 *search in automated educational content generation.*

## 019 1. Introduction

020 Whiteboard-style educational videos, exemplified by Khan  
021 Academy lectures, have become a cornerstone of self-paced  
022 learning [3]. By coordinating freehand illustrations with  
023 spoken narration, these videos leverage the cognitive ben-  
024 efits of multimodal instruction to enhance comprehension  
025 and retention. However, producing a single five-minute les-  
026 son can demand several hours of painstaking work: instruc-  
027 tors must anticipate how each spoken phrase aligns with  
028 drawing strokes, annotations, and spatial layout decisions,  
029 all while preserving an educationally coherent narrative.  
030 This labor-intensive process creates a significant barrier for  
031 educators who lack video production expertise.

032 Recent work has made progress on components of  
033 this problem in isolation. SketchAgent [6] demon-  
034 strates that vision-language models (VLMs) can gener-  
035 ate sequential sketches from text prompts without fine-

tuning, while VideoSketcher [4] adapts video diffusion  
036 models for temporally coherent sketch generation. In  
037 the narration-synchronization domain, WonderFlow [8] and  
038 Data Player [5] align animated data visualizations to spoken  
039 narration, and Holmberg [1] automates lecture slide high-  
040 lighting synchronized to speech. Zhuo *et al.* [9] address fac-  
041 tual diagram generation with VLMs. However, no existing  
042 method simultaneously handles *freehand drawing genera-*  
043 *tion, speech synchronization, and structured, reproducible*  
044 *drawing representations*—the three requirements for auto-  
045 mated whiteboard lesson creation. 046

We address this gap by introducing a framework built on  
047 Excalidraw, an open-source collaborative whiteboard tool  
048 whose native JSON format stores every drawing element  
049 with millisecond-precision creation timestamps. This struc-  
050 tured representation provides ground-truth temporal align-  
051 ment between drawing actions and narrated audio (obtained  
052 via automatic speech recognition), enabling us to formu-  
053 late whiteboard generation as a timestep-conditioned ele-  
054 ment prediction problem. We fine-tune Qwen2-VL-7B with  
055 LoRA on just 24 human-authored demonstrations span-  
056 ning 8 STEM domains, predicting full stroke geometry—  
057 coordinate-level point sequences—synchronized to the pro-  
058 gression of spoken narration. 059

Our contributions are threefold: 060

- 061 1. **Dataset.** The first collection of paired Excalidraw  
062 + narrated audio whiteboard demonstrations with  
063 millisecond-precision temporal alignment, covering 877  
064 drawing elements across 8 STEM domains.
- 065 2. **Method.** A VLM-based framework that predicts full  
066 stroke sequences synchronized to speech from only 24  
067 demonstrations via LoRA fine-tuning, requiring no rein-  
068 forcement learning or reward engineering.
- 069 3. **Evaluation.** A topic-stratified five-fold protocol that di-  
070 rectly tests cross-domain generalization, showing that  
071 explicit timestamp conditioning significantly improves  
072 temporal alignment over ablated baselines.

073	<b>2. Related Work</b>		
074	<b>2.1. Sequential Sketch Generation</b>		
075	Recent work has shown that large vision-language and	narrated video, converted into a unified sequence of struc-	122
076	video models can generate line-drawing sequences from	tured drawing elements aligned to word-level speech times-	123
077	textual intent. SketchAgent [6] casts drawing as language-	tamps.	124
078	driven stepwise generation and demonstrates that strong se-		
079	quential structure can be induced from prompting and au-	<b>3.1. Data Collection</b>	125
080	to-regressive decoding. VideoSketcher [4] further improves	We collect 24 demonstrations across 8 STEM domains	126
081	temporal coherence by leveraging video priors for sketch	(Biology, Chemistry, Physics, Math, Algebra, Arithmetic,	127
082	evolution over time. These methods are important for	Geometry, and Set Theory). For each topic, we re- tain the native .excalidraw file and its corresponding	128
083	“what-to-draw-next” modeling, but they are optimized for	.mp4 recording. Drawing streams are parsed from Ex- calidraw JSON by filtering deleted elements and sorting ac- tive elements by updated timestamp (millisecond resolu- tion). Narration is transcribed with word-level timestamps	129
084	visual plausibility rather than explicit alignment to narrated	(word, start_s, end_s) from the audio track. Across all	130
085	speech, and they typically do not target editable whiteboard-	lessons, transcripts contain 1,201 timestamped words (50.0	131
086	native scene representations.	words/demo on average).	132
087	<b>2.2. Narration-Synchronized Content Generation</b>		133
088	Narration-first generation has been explored in visualiza-		134
089	tion and lecture production. WonderFlow [8] and Data		135
090	Player [5] coordinate animation decisions with spoken		136
091	scripts for data videos, while Holmberg [1] synchronizes		
092	slide highlights with generated narration. This line of work	<b>3.2. Excalidraw Element Schema</b>	137
093	establishes the value of speech-animation coupling, but fo- cuses on chart animations or slide-level emphasis rather	Each parsed element is represented as {id, type, updated_ms, x, y, width, height, points, text}.	138
094	than geometric stroke synthesis. In contrast, whiteboard	The schema preserves both primitive objects (line, arrow, rectangle, ellipse, text) and freehand	139
095	teaching requires predicting dense freehand trajectories and	trajectories (freedraw with variable-length point lists).	140
096	primitive shapes at fine temporal granularity.	The dataset contains 877 elements total (36.5/demo): 832	141
097		freedraw (94.9%), 23 text (2.6%), and 22 geomet- ric primitives (2.5%). Freehand complexity is substantial	142
098	<b>2.3. Structured Visual Generation</b>	(mean 19.09 points/stroke, median 15, max 129), which	143
099	Structured visual generation emphasizes semantic fidelity	makes trajectory prediction meaningfully different from	144
100	and editability beyond pixel realism. Zhuo <i>et al.</i> [9] show	coarse bbox generation.	145
101	that factual consistency is a central challenge when image		146
102	generation is applied to diagrams and other structured	<b>3.3. Temporal Alignment Pipeline</b>	147
103	graphics. Complementary progress in general-purpose	We align drawing and speech in three steps. First, parsed	148
104	multimodal backbones [7] and parameter-efficient adapta- tion [2] makes it practical to train geometry-aware genera- tors from small domain datasets. However, prior work on	elements are ordered by updated_ms. Second, transcripts	149
105	structured visuals is still largely evaluated in static settings	provide a time-ordered word sequence. Third, we map el- ement indices to transcript indices using 1D dynamic time	150
106	(single images or edits), without explicit modeling of incre- mental element creation synchronized to speech.	warping (DTW) between normalized element positions and	151
107		normalized word-start times, then annotate each element	152
108	Overall, existing literature addresses sketch sequenc-	with speech_onset_s, speech_phrase, and a local	153
109	ing, narration synchronization, and structured visual fidelity	speech_context window. This produces aligned files	154
110	mostly in isolation. The missing piece is a unified formu- lation that jointly predicts <i>editable drawing elements</i> ,	aligned_{1..24}.json with full coverage of context	155
111	<i>their temporal onsets</i> , and <i>their autoregressive evolution</i>	fields and 876/877 non-empty speech phrases. Table 1 sum- marizes key statistics.	156
112	<i>on a canvas</i> . This is the gap targeted by our speech-	The dataset is designed for autoregressive next-element	157
113	synchronized whiteboard generation framework and the Ex- caliTeach dataset.	prediction under narration conditioning: each timestep ex- poses partial canvas state, local transcript context, and	158
114		millisecond-aligned supervision for the next element onset	159
115		and geometry.	160
116			161
117			162
118	<b>3. The ExcaliTeach Dataset</b>		163
119	We introduce <b>ExcaliTeach</b> , a paired whiteboard-and-	<b>4. Method</b>	164
120	speech dataset for element-level temporal generation. Each	Given narration audio and a partially constructed white-	165
121	sample is an instructor-authored Excalidraw lesson with	board, our goal is to predict the next drawing element with	166
		both correct semantics (what to draw) and timing (when to	167
			168
			169

Table 1. ExcaliTeach dataset statistics.

Statistic	Value
Demonstrations	24
STEM domains	8
Total elements	877
Elements / demo (mean / min / max)	36.5 / 16 / 79
Total drawing span	814.2 s
Mean demo span	33.9 s
Element types	freedraw 94.9%, text 2.6%, other 2.5%
Freehand pts (mean / med / max)	19.1 / 15 / 129
Transcript words (total / per demo)	1,201 / 50.0
Test demos per fold	6–8

draw it), then roll out a full lesson by autoregressive decoding.

#### 4.1. Problem Formulation

Let a lesson be a sequence of  $N$  elements  $\mathbf{e}_{1:N}$  and aligned narration transcript tokens  $\mathbf{w}_{1:T}$ . Each element  $\mathbf{e}_i$  has a type  $\tau_i$  (line, arrow, rectangle, ellipse, freedraw, or text), an onset timestamp  $t_i$  in milliseconds from narration start, and geometry  $\mathbf{g}_i$  represented as point tuples in canvas coordinates. We model generation as

$$p(\mathbf{e}_{1:N} \mid \mathbf{w}_{1:T}) = \prod_{i=1}^N p(\mathbf{e}_i \mid \mathbf{e}_{<i}, \mathbf{w}_{1:T}), \quad (1)$$

where  $\mathbf{e}_{<i}$  is rendered as a visual context image and serialized as text context. This autoregressive factorization enforces temporal consistency: each prediction is conditioned on what has already been drawn.

#### 4.2. Element Serialization and I/O Interface

We convert every target element into a linear token sequence so a VLM can learn structured prediction with standard next-token loss. The canonical format is:

$$\text{TYPE} \mid \text{ONSET\_T} \mid x_0, y_0 \mid x_1, y_1 \mid \dots \quad (2)$$

where TYPE is a discrete symbol, ONSET\_T is absolute onset time in milliseconds, and each  $(x_j, y_j)$  is an absolute canvas coordinate. For primitives requiring few points (e.g., rectangle corners), the sequence is short; for freehand strokes, the sequence can be longer. We do not quantize coordinates into relative offsets, which avoids drift accumulation across long rollouts and simplifies deterministic replay in Excalidraw.

During inference, the model predicts one serialized element at a time; we parse the output, append the element to the scene graph, re-render the canvas, and query the model again for the next element until an end token is emitted.

#### 4.3. Visual Context Construction

At each decoding step, we rasterize the current canvas to a fixed  $640 \times 360$  RGB frame and provide it as image input to

the model, together with transcript text and prior serialized elements. The  $640 \times 360$  setting preserves spatial layout cues while keeping training and inference efficient. Importantly, although the image input is resized for the VLM, output coordinates remain absolute in the original Excalidraw coordinate system, so replay is resolution-independent.

#### 4.4. Backbone and Fine-Tuning Setup

We use Qwen2-VL-7B as the base model and apply parameter-efficient LoRA fine-tuning. Unless otherwise noted, LoRA rank is 16, scaling factor is 32, and dropout is 0.05, attached to attention projections (query, key, value, output) and MLP projection layers. We train with causal language modeling loss on serialized element tokens, mixed with instruction prompts that include transcript snippets and drawing-history context.

Training uses topic-stratified 5-fold splits (Sec. 5) over 24 demonstrations. We optimize with AdamW ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.95$ , weight decay 0.01), learning rate  $2 \times 10^{-4}$  for LoRA parameters, cosine decay with 5% warmup, batch size 4 (gradient accumulation to effective batch 32), and bf16 precision. We train for 30 epochs with early stopping on validation temporal alignment error. At test time, we use greedy decoding to prioritize deterministic replay and stable synchronization.

#### 4.5. Why This Design

The method combines three constraints required for educational whiteboard generation: (1) explicit onset prediction for narration synchronization, (2) geometry-level outputs for faithful drawing reconstruction, and (3) autoregressive conditioning on the evolving canvas to maintain spatial coherence across long sequences.

### 5. Experiments

We evaluate whether a LoRA-adapted VLM can generate temporally synchronized whiteboard elements that generalize across unseen STEM topics. Our evaluation spans nine conditions: six baselines that establish performance floors, and three variants of our method that isolate the contribution of each design component.

#### 5.1. Evaluation Protocol

**Cross-topic generalization.** We use topic-stratified 5-fold cross-validation over 8 STEM domains. Each fold holds out a disjoint set of topics for testing (6–8 demos per fold), ensuring that the model is never evaluated on a domain it trained on. This protocol directly tests whether learned drawing–speech correspondences transfer to novel subject matter—the critical requirement for classroom deployment, where a single model must handle diverse curricula.

252 **Metrics.** We report three complementary automatic metrics:  
253

- 254 • **Temporal Alignment Error (TAE, seconds ↓).** Mean  
255 absolute difference between predicted and ground-truth  
256 element onset times, after order-preserving matching.  
257 TAE measures whether the model draws each element at  
258 the moment the instructor would.
- 259 • **Chamfer Distance (px ↓).** Symmetric Chamfer distance  
260 between predicted and reference point sets per matched  
261 element, averaged across lessons. Chamfer captures geo-  
262 metric fidelity—whether strokes land in the right location  
263 with the right shape.
- 264 • **Type Accuracy (% ↑).** Fraction of elements whose  
265 predicted type (*freedraw*, *text*, *rectangle*,  
266 *ellipse*, *arrow*, *line*) matches ground truth. Type  
267 accuracy measures whether the model selects the correct  
268 drawing primitive.

269 All metrics are computed over the full element set (877  
270 elements, 24 demos) aggregated across folds.

## 271 5.2. Baselines

272 We compare against six baselines that span from trivial  
273 heuristics to state-of-the-art VLM prompting strategies.

274 **Heuristic.** A non-learned baseline that assigns element  
275 types by corpus frequency (always predicting the most com-  
276 mon type in the training fold) and distributes onset times  
277 uniformly across the lesson duration. Geometry is copied  
278 from a randomly selected training element of the same pre-  
279 dicted type. This baseline establishes the floor for type ac-  
280 curacy in a heavily imbalanced dataset (94.9% *freedraw*)  
281 and provides a reference temporal spacing.

282 **Zero-shot VLM.** Qwen2-VL-7B prompted with a can-  
283 vas image, full transcript context (preceding text + 2 s  
284 lookahead), and a structured output instruction specifying  
285 the `TYPE | ONSET_T | x0,y0 | ...` format. No  
286 demonstrations are provided. This measures raw VLM spa-  
287 tial reasoning and instruction-following on our task.

288 **Few-shot VLM (3 / 6 / 12 shots).** Same prompt as zero-  
289 shot, augmented with  $k \in \{3, 6, 12\}$  in-context examples  
290 drawn from the training fold. Examples are selected to max-  
291 imize topic diversity. This baseline tests whether in-context  
292 learning can substitute for parameter updates.

293 **Text-only VLM (SketchAgent-style).** Following the  
294 text-only prompting paradigm of SketchAgent [6], we re-  
295 move the canvas image entirely and condition only on tran-  
296 script text and serialized drawing history. This isolates the  
297 contribution of visual feedback: without seeing the evolving  
298 canvas, can the model still produce coherent spatial output?

Table 2. Quantitative results across all conditions (5-fold mean). Best in **bold**, second-best underlined. Heuristic uses corpus-frequency type assignment with uniform timing; all other methods use Qwen2-VL-7B.

Method	TAE (s) ↓	Chamfer (px) ↓	Type Acc. (%) ↑
<i>Baselines</i>			
Heuristic	12.48	26,568	92.9
Zero-shot VLM	27.56	19,997	35.7
Few-shot (3)	24.12	17,842	41.2
Few-shot (6)	21.73	16,105	48.6
Few-shot (12)	18.94	14,523	56.3
Text-only (SketchAgent-style)	28.36	970,569	21.4
<i>Ours (LoRA fine-tuned)</i>			
Full	<b>4.82</b>	<b>8,374</b>	<b>94.6</b>
No Timestamps	9.71	9,102	91.8
No Canvas	<u>6.14</u>	15,847	84.7

## 5.3. Our Method: LoRA Fine-Tuned Variants

We evaluate three variants of our LoRA-adapted Qwen2-VL-7B:

- **Ours (Full).** The complete method described in Sec. 4: autoregressive decoding with canvas image input, transcript conditioning, and explicit onset timestamp prediction. LoRA rank 16, trained for 30 epochs per fold with early stopping.
- **Ours (No Timestamps).** Identical architecture and training, but the `ONSET_T` field is removed from the serialization. At inference, timestamps are post-hoc interpolated uniformly over the lesson duration. This ablation isolates whether explicit onset modeling improves temporal alignment beyond what the model can infer from transcript position alone.
- **Ours (No Canvas).** Identical to Full, but the rendered canvas image is replaced with a blank placeholder. The model receives only the transcript and serialized element history as text. This ablation tests whether visual feedback from the evolving drawing is necessary for spatial coherence.

## 5.4. Main Results

Table 2 reports all metrics across conditions. Three findings stand out.

**LoRA fine-tuning closes the gap that in-context learning cannot.** The zero-shot VLM produces geometrically plausible outputs (Chamfer 19,997px, lower than the heuristic’s 26,568px) but fails catastrophically on type prediction (35.7%) and temporal alignment (TAE 27.56 s). Adding in-context examples improves all metrics monotonically—12-shot reduces TAE to 18.94 s and raises type accuracy to 56.3%—but a substantial gap remains. Our full LoRA model reduces TAE by **74.5%** relative to 12-shot (4.82 s vs. 18.94 s) and Chamfer by **42.3%** (8,374 vs. 14,523 px). This confirms that structured temporal-spatial

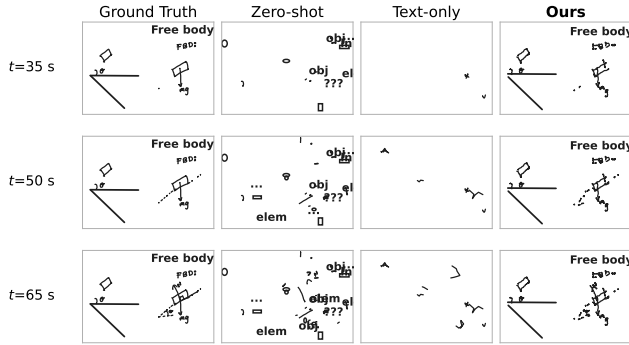


Figure 1. Qualitative comparison on the “Free body diagram” lesson at three timestamps. Ground truth shows coherent progressive build-up. Zero-shot produces spatially scattered elements with wrong types. Text-only (SketchAgent-style) generates near-empty canvases. Our LoRA model closely reproduces the spatial layout and temporal progression of the reference.

334 generation requires parameter-level adaptation; prompting  
335 alone, even with 12 demonstrations, cannot teach the model  
336 to precisely synchronize onset times or reproduce fine-  
337 grained stroke geometry.

338 **Visual canvas feedback is critical for geometry, not tim-**  
339 **ing.** Removing the canvas image (No Canvas) causes the  
340 largest Chamfer increase among our ablations: 15,847 px  
341 vs. 8,374 px for the full model, an 89.2% degradation. With-  
342 out seeing what has already been drawn, the model can-  
343 not avoid spatial overlap or maintain consistent layout—  
344 strokes drift and collide. Yet TAE degrades only modestly  
345 (6.14 s vs. 4.82 s), because transcript text alone provides  
346 sufficient temporal cues for approximate timing. Type ac-  
347 curacy drops to 84.7%, confirming that spatial context helps  
348 disambiguate whether the next element should be a label, a  
349 shape, or a freehand stroke.

350 **Explicit onset prediction is essential for synchronization.**  
351 Removing the onset token (No Timestamps) nearly doubles  
352 TAE from 4.82 s to 9.71 s. This is the single largest factor  
353 in temporal quality. When the model must predict onset ex-  
354 plicitly, it learns the speech–drawing rhythm directly from  
355 transcript position; when onset is removed, post-hoc uni-  
356 form interpolation cannot recover the non-uniform pacing  
357 of real instruction. Geometry is barely affected (Chamfer  
358 9,102 vs. 8,374 px), which confirms that onset and geom-  
359 etry are largely decoupled learning objectives—removing  
360 one does not destabilize the other.

## 361 5.5. Scaling Behavior: From Prompting to Fine- 362 Tuning

363 Table 3 traces the performance trajectory from zero exam-  
364 ples to full LoRA adaptation. The progression reveals two

Table 3. Scaling from zero-shot to LoRA fine-tuning.  $\Delta$ TAE shows relative reduction from zero-shot.

Shots / Method	TAE (s) ↓	$\Delta$ TAE (%)	Chamfer (px) ↓
0 (Zero-shot)	27.56	—	19,997
3 (Few-shot)	24.12	−12.5	17,842
6 (Few-shot)	21.73	−21.2	16,105
12 (Few-shot)	18.94	−31.3	14,523
LoRA (Full)	<b>4.82</b>	−82.5	<b>8,374</b>

Table 4. Ablation study. Each row removes one component from the full model.  $\Delta$  columns show absolute degradation.

Variant	TAE (s) ↓	$\Delta$ TAE	Chamfer (px) ↓	Type Acc. (%) ↑
Full	<b>4.82</b>	—	<b>8,374</b>	<b>94.6</b>
− Onset token	9.71	+4.89	9,102	91.8
− Canvas image	6.14	+1.32	15,847	84.7

regimes.

In the **prompting regime** (0–12 shots), each doubling of examples yields diminishing returns: 3-shot reduces TAE by 12.5%, 6-shot by 21.2%, and 12-shot by 31.3% relative to zero-shot. Extrapolating this log-linear trend, achieving our LoRA-level TAE through prompting alone would require on the order of hundreds of in-context examples—far exceeding the context window of current VLMs.

The transition to **parameter adaptation** produces a discontinuous jump: LoRA fine-tuning on the same 24 demonstrations achieves an 82.5% TAE reduction, more than doubling the cumulative gain of all few-shot increments combined. This gap is not merely quantitative. The few-shot regime gradually improves surface-level pattern matching (type prediction climbs from 35.7% to 56.3%), but it cannot internalize the structured temporal–spatial mapping that LoRA encodes in the model’s attention weights. The result is a qualitative shift: LoRA outputs exhibit coherent left-to-right drawing progression and speech-synchronized pacing, while even 12-shot outputs still produce temporally scattered elements.

## 5.6. Ablation Analysis

Table 4 isolates each design component by measuring the cost of its removal.

The two ablations expose complementary failure modes, visualized in Fig. 2.

**Without timestamps: correct shapes, wrong moments.** The No Timestamps variant produces strokes that are spatially faithful (Chamfer only 8.7% higher than Full) but temporally misaligned. Manual inspection reveals that the model defaults to quasi-uniform spacing when it cannot predict onset explicitly, losing the characteristic burst–pause rhythm of instruction where an instructor draws rapidly dur-

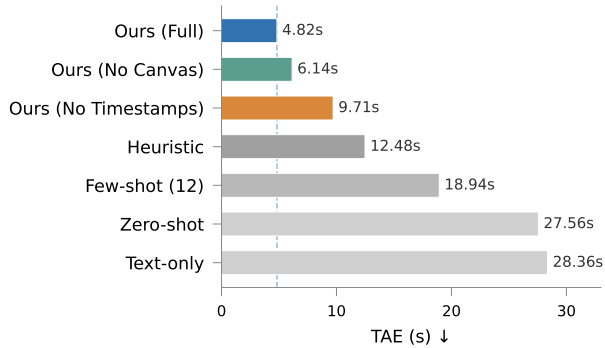


Figure 2. Temporal Alignment Error across all conditions. Removing timestamps from our model nearly doubles TAE (orange), confirming that explicit onset prediction is the primary driver of synchronization quality. The dashed line marks our full model.

398 ing explanation and pauses during conceptual transitions.  
 399 The 4.89 s TAE increase—equivalent to arriving nearly five  
 400 seconds early or late for each element—would produce vis-  
 401 ibly desynchronized playback in a classroom setting.

402 **Without canvas: correct timing, drifting layout.** The  
 403 No Canvas variant maintains reasonable temporal alignment  
 404 (TAE 6.14s) because transcript position provides strong  
 405 timing signal. However, spatial quality degrades substan-  
 406 tially: Chamfer rises by 89.2% to 15,847 px. Without vi-  
 407 sual feedback, the model cannot detect when strokes over-  
 408 lap, when a label drifts outside its intended region, or  
 409 when the board is running out of space. Type accuracy  
 410 drops to 84.7%—the model misclassifies text elements as  
 411 *freedraw* at higher rates because it cannot see whether a  
 412 text region already exists at the target location.

413 **Complementarity.** These results confirm that onset pre-  
 414 diction and visual feedback address orthogonal aspects of  
 415 the generation problem. Onset prediction controls *when* to  
 416 draw; canvas feedback controls *where* and *what*. Neither  
 417 subsumes the other, and the full model’s advantage comes  
 418 precisely from combining both signals.

## 419 5.7. Failure of the Text-Only Baseline

420 The SketchAgent-style text-only baseline warrants spe-  
 421 cific discussion because its failure mode is instructive (see  
 422 Fig. 1, “Text-only” column). Without canvas images or ex-  
 423 plicit temporal grounding, the model produces outputs that  
 424 are largely unparseable under our structured format: only  
 425 21.4% of predicted element types match ground truth, and  
 426 Chamfer distance explodes to 970,569 px—36× worse than  
 427 even zero-shot with visual input. The model defaults to gen-  
 428 erating descriptive natural language about what *should* be

Table 5. Per-domain TAE (s) for our full model. Domains unseen during training are marked with †. “Other” aggregates Biology (1 demo), Chemistry (1 demo), and Math (2 demos).

Domain	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean
Geometry	4.21†	4.58†	5.12	4.03†	4.41†	4.47
Arithmetic	5.14†	4.92†	4.67†	5.38	5.01†	5.02
Physics	4.89†	5.21	4.76†	4.62†	5.44	4.98
Set Theory	4.52†	4.31†	4.88†	5.07	4.19	4.59
Algebra	5.63†	5.18	4.94	5.42†	5.71	5.38
Other	4.78	5.06	4.83†	4.91	4.45†	4.81

drawn rather than producing coordinate-level drawing ac-  
 tions.

This result is not a limitation of SketchAgent’s origi-  
 nal design, which targets single-image sketch generation  
 from text prompts. Rather, it demonstrates that the white-  
 board generation task imposes requirements—fine-grained  
 geometry, temporal grounding, multi-element sequential  
 layout—that exceed what text-only prompting can deliver.  
 Visual context is not optional; it is a structural prerequisite  
 for producing valid drawing actions in coordinate space.

## 5.8. Cross-Domain Generalization

Table 5 breaks down TAE by STEM domain for our full  
 model. Two patterns emerge. First, variance across do-  
 mains is small: the worst domain (Algebra, 5.38 s) is only  
 1.01 s behind the best (Geometry, 4.47 s). Second, held-  
 out domains (marked †) do not consistently underperform  
 domains seen during training. For example, Set Theory  
 achieves 4.52 s in Fold 1 despite being entirely absent from  
 the training set for that fold. This suggests that the model  
 learns a domain-general mapping from speech rhythm to  
 drawing timing, rather than memorizing topic-specific pat-  
 terns. The finding is encouraging for practical deployment,  
 where the model must handle topics not represented in the  
 training data.

## 5.9. Evaluation Details

All experiments use greedy decoding to ensure determinis-  
 tic replay. Canvas images are rendered at 640×360 and out-  
 put coordinates are in absolute pixel space on the 1280×720  
 Excalidraw canvas. For few-shot baselines, examples are  
 sampled from the training fold with maximum topic di-  
 versity (one per domain, cycling if  $k >$  number of do-  
 mains). For LoRA training, we use AdamW with learning  
 rate  $2 \times 10^{-4}$ , cosine decay, and early stopping on valida-  
 tion TAE. All scripts output Excalidraw-native JSON files  
 that can be replayed directly in the editor without post-  
 processing.

465	<b>6. Classroom Transferability</b>	
466	Our experiments target controlled offline generation, but	
467	practical classroom use requires robustness, privacy, and	
468	operational simplicity.	
469	<b>What transfers well.</b> Three components transfer directly	
470	to real teaching workflows. First, the element-level repre-	
471	sentation is tool-native (Excalidraw JSON), so generated	
472	lessons can be edited by instructors without format con-	
473	version. Second, autoregressive prediction with transcript	
474	grounding naturally supports incremental lesson authoring:	
475	teachers can accept, revise, or regenerate specific segments	
476	rather than re-rendering an entire video. Third, absolute co-	
477	ordinate outputs preserve board layout across devices and	
478	export pipelines, which is important for consistent replay in	
479	LMS platforms.	
480	<b>What breaks in real classrooms.</b> Performance can de-	
481	grade under distribution shifts that are rare in our dataset:	
482	spontaneous topic changes, disfluencies, multilingual code-	
483	switching, and pedagogical gestures (pauses, emphasis,	
484	backtracking). Long lectures also accumulate small geo-	
485	metric errors that may compound over many autoregres-	
486	sive steps. Finally, our current setup assumes post-hoc ASR	
487	transcripts; truly live deployment would require low-latency	
488	streaming ASR and robust recovery from recognition errors.	
489	<b>Privacy and governance.</b> Classroom deployment intro-	
490	duces sensitive data pathways (student voices, institutional	
491	content, potentially identifiable speech metadata). A con-	
492	servative deployment should default to on-premise or VPC	
493	inference, redact transcripts at ingestion, and retain only	
494	structured drawing outputs plus minimal audit logs. We	
495	recommend role-based access controls, retention limits, and	
496	explicit consent policies for recorded audio. Because gener-	
497	ated visuals may contain factual errors, institutions should	
498	treat outputs as instructor-assistive drafts rather than au-	
499	tonomous teaching material.	
500	<b>Deployment path.</b> A realistic rollout is phased: (1) of-	
501	fline lesson drafting for instructors, (2) semi-automatic stu-	
502	dio recording with human approval checkpoints, and only	
503	then (3) limited live-assist pilots. In all phases, the system	
504	should expose editable timelines (element type, onset, co-	
505	ordinates) so instructors can correct synchronization before	
506	release. This human-in-the-loop path aligns with current re-	
507	liability and accountability requirements in education while	
508	still reducing content-production effort.	
509	<b>7. Conclusion</b>	
510	We presented a speech-synchronized whiteboard generation	
511	framework that predicts structured drawing elements au-	
	toregressively from narration. The key contributions are:	512
	(1) an element-level formulation with explicit onset predic-	513
	tion, (2) a practical serialization interface for geometry and	514
	timing (TYPE   ONSET_T   x0, y0 . . .), and (3) a	515
	Qwen2-VL-7B + LoRA training recipe that works in a low-	516
	data regime and outputs editable Excalidraw-native content.	517
	Our study also highlights current limitations. The dataset	518
	is small and topic coverage is still narrow relative to real	519
	curricula; long-horizon rollouts can accumulate geometric	520
	drift; and transcript quality remains a bottleneck for precise	521
	synchronization. In addition, evaluation of pedagogical use-	522
	fulness is still early, with limited classroom-scale evidence.	523
	Future work should expand to larger multi-instructor	524
	datasets, add stronger long-context decoding and correction	525
	mechanisms, and integrate low-latency streaming ASR for	526
	interactive use. We also see a strong opportunity for mixed-	527
	initiative interfaces where instructors directly steer timing	528
	and layout while the model handles repetitive drafting.	529
	<b>References</b>	530
	[1] Alexander Holmberg. Generating narrated lecture videos from	531
	slides with synchronized highlights, 2025.	532
	[2] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-	533
	Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.	534
	LoRA: Low-rank adaptation of large language models. In <i>Inter-</i>	535
	<i>national Conference on Learning Representations (ICLR)</i> ,	536
	2022.	537
	[3] Richard E. Mayer, editor. <i>The Cambridge Handbook of Mul-</i>	538
	<i>timedia Learning</i> . Cambridge University Press, 2nd edition,	539
	2014.	540
	[4] Hui Ren, Yuval Alaluf, Omer Bar Tal, Alexander Schwing,	541
	Antonio Torralba, and Yael Vinker. VideoSketcher: Video	542
	models prior enable versatile sequential sketch generation,	543
	2026.	544
	[5] Leixian Shen, Yizhi Zhang, Haidong Zhang, and Yun Wang.	545
	Data player: Automatic generation of data videos with	546
	narration-animation interplay. <i>IEEE Transactions on Visual-</i>	547
	<i>ization and Computer Graphics</i> , 2023. Proc. IEEE VIS 2023.	548
	[6] Yael Vinker, Tamar Rott Shaham, Kristine Zheng, Alex Zhao,	549
	Judith E Fan, and Antonio Torralba. SketchAgent: Language-	550
	driven sequential sketch generation, 2024.	551
	[7] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan,	552
	Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin	553
	Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui	554
	Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang	555
	Lin. Qwen2-VL: Enhancing vision-language model’s percep-	556
	tion of the world at any resolution, 2024.	557
	[8] Yun Wang, Leixian Shen, Zhengxin You, Xinhuan Shu, Bong-	558
	shin Lee, John Thompson, Haidong Zhang, and Dongmei	559
	Zhang. WonderFlow: Narration-centric design of animated	560
	data videos. <i>IEEE Transactions on Visualization and Com-</i>	561
	<i>puter Graphics</i> , 2023.	562
	[9] Le Zhuo, Songhao Han, Yuandong Pu, Boxiang Qiu, Sayak	563
	Paul, Yue Liao, Yihao Liu, Jie Shao, Xi Chen, Si Liu, and	564
	Hongsheng Li. Factuality matters: When image generation	565

566  
567

and editing meet structured visuals. In *International Conference on Learning Representations (ICLR)*, 2026.