
Generative Timelines for Instructed Visual Assembly

Alejandro Pardo^{1,2}, Jui-Hsien Wang², Bernard Ghanem¹, Josef Sivic^{2,3}, Bryan Russell²
Fabian Caba Heilbron²

¹AI Initiative, KAUST

²Adobe Research

³CIIRC CTU*

<https://sites.google.com/kaust.edu.sa/timeline-assembler>

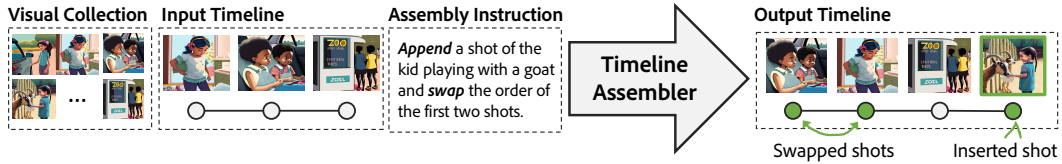


Figure 1: **Instructed Visual Assembly.** Given a visual collection, an input timeline, and an assembly instruction, our model (called the Timeline Assembler) performs the instructed assembly task and generates an output timeline with the desired edits. The collection comprises various media elements, such as video clips or images. The timeline is a sequential arrangement of these elements.

Abstract

The objective of this work is to manipulate visual timelines (*e.g.*, a video) through natural language instructions, making complex timeline editing tasks accessible to non-expert or potentially even disabled users. We call this task *Instructed visual assembly*. This task is challenging as it requires (i) identifying relevant visual content in the input timeline as well as retrieving relevant visual content in a given input (video) collection, (ii) understanding the input natural language instruction, and (iii) performing the desired edits of the input visual timeline to produce an output timeline. To address these challenges, we propose the Timeline Assembler, a generative model trained to perform instructed visual assembly tasks. The contributions of this work are three-fold. First, we develop a large multimodal language model, which is designed to process visual content, compactly represent timelines and accurately interpret timeline editing instructions. Second, we introduce a novel method for automatically generating datasets for visual assembly tasks, enabling efficient training of our model without the need for human-labeled data. Third, we validate our approach by creating two novel datasets for image and video assembly, demonstrating that the Timeline Assembler substantially outperforms established baseline models, including the recent GPT-4o, in accurately executing complex assembly instructions across various real-world inspired scenarios.

*Czech Institute of Informatics, Robotics and Cybernetics at the Czech Technical University in Prague.

¹Work partially done during an internship at Adobe.

1 Introduction

Imagine returning from a zoo trip and finding a visual timeline, *e.g.*, in the form of a short video, of highlights automatically generated by your device [32, 9]. While reviewing the timeline, you notice some shots are misplaced and a memorable moment with that kid petting a goat is missing. Traditionally, making such modifications would require finding the additional shot in your video collection as well as potential cumbersome interactions with video editing tools. This may be even harder for users with small-screen devices or users with disabilities for whom it is hard to interact with traditional interfaces. A more intuitive approach, illustrated in Figure 1, involves using natural language to direct edits. The users simply state their desired changes, much like setting a calendar event by voice. Such a system uses the collection of visual assets, the current timeline, and the provided instruction to return the refined timeline with the requested edits.

We define this capability as instructed visual assembly, a process that involves the automated editing of a visual timeline in response to user-provided natural language instructions. To effectively automate this task, the system requires a comprehensive understanding of three key elements: the instruction itself, the existing visual timeline, and the collection of visual assets. For instance in Figure 1, if a user instructs the system to *"swap the first two clips and add the shot of the kid with the goat,"* the system must first understand the multimodal context by interpreting the wording of the instruction and relating it to the visual data in the timeline and the collection. The system must then identify the specific elements to be edited, such as the initial two shots, and retrieve the additional clip of the child with the goat from the collection. This process involves deep multimodal understanding to determine not only what the user intends but also how these intentions translate into direct manipulation of visual content. Despite these challenges, instructed visual assembly can enable intuitive video creation interfaces for novices, and users with disabilities, mitigating the complexities associated with mastering traditional video editing tools and interfaces, as well as easing the management of visual content. Moreover, it can ease the editing of videos on devices with small screens where simple operations like drag and drop become cumbersome.

This is the first work that explicitly addresses the task of instructed visual assembly. While this task shares similarities with language-based video editing [12, 37], which consists of editing the pixels of a source video based on language instructions, instructed visual assembly distinctively focuses on executing instructions to arrange visual elements in a timeline. There are other related tasks in the video assembly space, such as automatic shot transitioning [36, 35] and B-roll recommendation [18, 50, 45]. However, these tasks primarily provide creative guidance, diverging from the instructed visual assembly goal of composing timelines through natural language instructions.

In this paper, we present a method for teaching a generative model, named Timeline Assembler, to follow assembly instructions and generate timelines with the appropriate edits. The Timeline Assembler builds on Large Language Models (LLMs) to leverage their remarkable skills in following instructions [49, 47] and interpreting multimodal content [28, 54]. Specifically, we adapt a multimodal LLM [54] to handle the nuances of instructed visual assembly. This adaptation presents two major challenges. Firstly, existing multimodal LLMs are designed to process a single image or video. Thus, devising a representation for handling visual collections and timelines using LLMs remains an open challenge. Secondly, as shown by our experiments, existing multimodal LLMs struggle to understand and execute visual assembly instructions, hinting that further tuning for the assembly task is required. However, we face the challenge of doing so without human-labeled data, which is hard to acquire at scale for such specialized tasks.

Our work makes the following contributions to address the challenges above:

- (1) We design a multimodal LLM architecture to process visual collections, encode timelines, and ingest natural language instructions to generate edited output timelines. Our architecture represents each image or video clip in the collection with a unique identifier and a visual representation compatible with the LLM’s input space [28]. Secondly, we encode the timeline using the unique identifiers of each arranged visual element. Finally, we map the LLM’s output tokens directly back to a timeline, associating each token to its relevant visual element. This design not only ensures a compact representation of the visual collection and the timeline, but also facilitates a straightforward reconstruction of timelines from the model’s output tokens.
- (2) We propose a new approach to automatically generate a paired dataset (input/output) for instructed visual assembly. Our method programmatically creates a collection, input/output timelines, and

assembly instructions from an input visual sequence and task candidates. Using this approach, we generate data to learn the projection layer and Low-Rank Adapters (LoRA) [15] for effective LLM performance on assembly tasks. Our training approach is not only human-labels-free but also efficient and enables a wide range of timeline editing tasks.

(3) We construct two datasets, one for image sequence assembly (of still images) and one for video assembly, to evaluate instructed visual assembly. On both datasets, our method substantially outperforms baseline approaches, including powerful LLMs such as GPT3.5, highlighting the benefits of our approach. Moreover, we show that our model can match or even exceed the performance of specialized (single-task) assembly models, perform equally well regardless of the length of the timelines, and execute multiple complex instructions at a time.

2 Related Work

Multimodal LLMs as Multi-task Interfaces. Our work shares a similar spirit with recent trends that use LLM capabilities to perform multiple tasks for different applications [14]. Several methods have been proposed for image and language tasks [54, 2, 28, 22, 8, 20], video understanding [31, 24, 7, 43, 29], multimodal understanding [51, 52], recommendation systems [13], and robotics and motion planning [19, 26, 4, 10, 16]. Our Timeline Assembler distinguishes itself by specializing in the domain of visual timelines, interpreting natural language to modify sequences of visual data. This unique application extends the use of LLMs beyond physical tasks to the manipulation of multimedia content.

Video Assembly. Several computer vision techniques have been proposed for automated video editing tasks [39, 36, 42, 3, 21, 35, 6, 40]. Likewise, multiple works have proposed to tackle video assembly. Li et al. [25] propose multi-shot vlog assembly, drawing parallels to our approach where an initial sequence guides subsequent shot selections from a candidate pool. Furthermore, other works [46, 48, 50, 30] have developed tools that assemble videos from input queries, leveraging multiple components for an intuitive video creation experience. While these works emphasize on providing creative guidance, our model introduces a user-guided approach to video assembly, enabling a novel interface to manipulate and arrange timelines with language instructions.

Instruction Fine-tuning and Aligning with User Intent. Our paper builds on methods that adapt language models to closely following human instructions [38, 5]. Wei *et al.* [49] first introduced instruction-tuning, enhancing the usability and multi-tasking of LLMs. [34] expanded this by integrating Reinforcement Learning from Human preferences (RLHF). [47] demonstrated that GPT3-generated instructions could achieve results similar to InstructGPT through self-instruct fine-tuning. Our work uses similar ideas and tailors instruction fine-tuning for visual assembly tasks. We do so by automatically gathering novel visual assembly instruction data that serves as training data to our LoRA [15] fine-tuned multimodal LLM.

3 Timeline Assembler

3.1 Timeline Assembler Architecture

Our goal is to design the Timeline Assembler to integrate visual and textual data to generate output timelines that precisely align with the given assembly instructions. The Timeline Assembler generates an edited timeline \tilde{S} from an assembly instruction q , an initial timeline S , and a collection of video clips C . With instructions often phrased in multiple ways, ensuring consistent output results poses a considerable challenge. To tackle these challenge, the Timeline Assembler leverages an instruction-following Large Language Model (LLM) that reasons over text and visual tokens. Our overall architecture, illustrated in Figure 2, comprises a set of tokenizers that convert multimodal inputs into a set of tokens \mathbf{X} . These tokens are then passed to an LLM, which returns an output sequence of tokens $\mathbf{X}_{\tilde{S}}$ that represent a timeline in the LLM space. Finally, we reconstruct the output timeline \tilde{S} from $\mathbf{X}_{\tilde{S}}$. We describe each of these components in detail next.

Collection Tokenizer. Our goal is to represent a collection of visual elements as an array of unique identifier tokens and visual tokens to enable an LLM to process the visual content in the collection. Each identifier token serves to distinctly identify the visual elements within the collection, while the

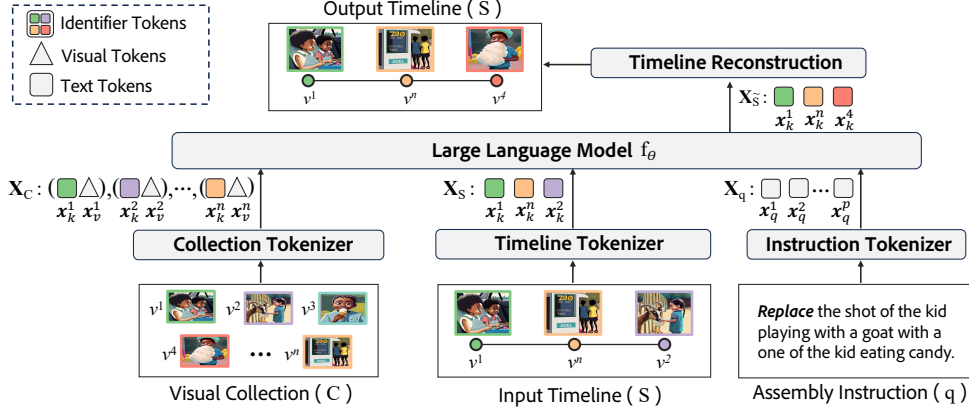


Figure 2: **Timeline Assembler Architecture.** We design a multimodal architecture to execute visual assembly instructions to generate visual timelines. Our model takes as inputs: a collection of images/videos C , a timeline S , and an assembly instruction q . Each image/video in the collection is represented with a unique identifier token x_k^i (color-coded) and a visual token x_v^i , forming the sequence \mathbf{X}_C . The input timeline is represented with the sequence of tokens \mathbf{X}_S , which comprises the list of identifier tokens of the images/videos in the timeline. The assembly instruction is tokenized into \mathbf{X}_q . Given the input tokens, the task of the Large Language Model is to generate output timeline tokens $\mathbf{X}_{\tilde{S}}$, which are reconstructed into the output timeline \tilde{S} .

visual tokens encapsulate the visual information. With this design, the LLM can reference the visual representations associated with each element in the collection.

Let $C = \{v^1, v^2, \dots, v^n\}$ denote the visual collection, where v^i represents a visual element composed of an image or video clip. The collection tokenizer has two key components: a mapping function \mathcal{H} that generates a unique identifier token x_k^i for each visual element v^i ; and a visual encoder that embeds every visual input v^i to produce a visual token x_v^i . The mapping function \mathcal{H} obtains a unique text token x_k^i by assigning a token from a set of previously tokenized integers. In practice, this mapping function operates as a look-up table that helps to assign (and find) the unique identifier token of a given visual element. The visual encoder ingests a visual element v^i and outputs a visual token x_v^i . In practice, a pretrained visual encoder g ingests and extracts visual representations from the input visual elements, and a projection layer $h_\gamma(\cdot)$ maps these visual representations into a visual token that is aligned with the input space of the LLM [33]. Each visual element v^i is then represented as a tuple (x_k^i, x_v^i) consisting of a unique identifier token and a visual token, respectively, such that the entire array of collection tokens $\mathbf{X}_C = [(x_k^1, x_v^1), (x_k^2, x_v^2), \dots, (x_k^n, x_v^n)]$.

Timeline Tokenizer. The goal of the timeline tokenizer is to map a sequence of visual elements within a timeline to their corresponding identifier tokens, allocated during the collection’s initial tokenization. By utilizing these pre-existing tokens, the tokenizer avoids the need for re-tokenizing visual elements as they appear in the timeline. Doing so prevents redundant visual tokenization, enabling the LLM to handle in practice more extensive collections and longer timelines. As a result, the visual elements (images or videos) are represented only once to save space. Their corresponding identifier token acts as a reference (or pointer) to the (typically high-dimensional) visual token.

In detail, let $S = \{v^{S1}, v^{S2}, \dots, v^{Sl}\}$ denote the input timeline, where v^{Si} denotes the i -th visual element in the timeline, and l is the length of the timeline. Using the mapping function \mathcal{H} , we can retrieve the identifier token for each visual element v^{Si} and construct the tokenized timeline $\mathbf{X}_S = [x_k^{S1}, x_k^{S2}, \dots, x_k^{Sl}]$.

Instruction Tokenizer. Given an input instruction text prompt q , the goal of the tokenizer is to map the sequence of strings/words into a discrete set of p tokens such that: $\mathbf{X}_q = [x_q^1, x_q^2, \dots, x_q^p]$, where \mathbf{X}_q is a sequence of tokens that represents the full text prompt q . In practice, we leverage the byte-level BPE text tokenizer [41] of our large language model.

Finally, the input to the LLM is the union of token sequences \mathbf{X} , defined as $\mathbf{X} = [\mathbf{X}_C, \mathbf{X}_S, \mathbf{X}_q]$.

Large Language Model. Our goal is to develop a model that can assemble and edit timelines of visual sequences. The model must be able to understand natural language instructions and a

multimodal context representing a visual collection and the input timeline. To do so, we leverage multimodal LLMs given their capabilities at managing multiple multimodal tasks [28, 8], reasoning over long sequences [7, 27], and encapsulating knowledge from a plethora of sources [34, 44].

Therefore, we employ a Large Language Model $f_\theta(\cdot)$, parameterized by θ , to generate an updated timeline sequence $\mathbf{X}_{\tilde{S}}$ in response to an input instruction. The model takes a sequence of previously defined multi-modal tokens \mathbf{X} as input and, at test time, outputs a sequence of identifier tokens representing the updated timeline’s visual elements: $f_\theta(\mathbf{X}) \rightarrow \mathbf{X}_{\tilde{S}} = [\mathbf{x}_k^{\tilde{S}1}, \mathbf{x}_k^{\tilde{S}2}, \dots, \mathbf{x}_k^{\tilde{S}Q}]$, where Q is the output timeline length, and $\mathbf{x}_k^{\tilde{S}i}$ is the identifier token for the i -th element in the updated timeline.

Timeline Reconstruction. Given the output tokens $\mathbf{X}_{\tilde{S}}$, produced by the LLM $f_\theta(\cdot)$, the goal of this step is to reconstruct the output timeline \tilde{S} . To do so, we map each output token, which in practice are identifier tokens, to their corresponding visual elements using a reverse mapping operation from \mathcal{H} .

3.2 Constructing Visual Assembly Tasks

Our goal is to train the Timeline Assembler to effectively perform various visual assembly tasks. To facilitate this, we define a suite of atomic operations that act as the foundational elements for constructing and manipulating visual timelines. These operations are:

- *Insert (in)*: Add an element from the collection to the timeline.
- *Remove (rm)*: Delete an entry from the timeline.
- *Replace (rp)*: Substitute one element in the timeline with another from the collection.
- *Swap (sw)*: Exchange positions of two elements within the timeline.

To reference elements, either in the timeline or the collection, we use two types of cues:

- *Positional Cue (p)*: Refers to an element by its identifier or position in the timeline.
- *Semantic Cue (r)*: Refers to an element through a language-based description of its visual content.

These operations and cues combine to form eight distinct assembly tasks $T = t_c$, where t represents the operations *in*, *rm*, *rp*, *sw* and c denotes the types of cues *p*, *r*. An illustration of each one of these tasks can be found in Figure ?? . To generate training data for these tasks, we apply a transformation function ϕ_{t_c} to an initial timeline S^i . This function manipulates the timeline to produce a modified version \tilde{S}^i . For example, to train the model for the "remove" operation, we artificially introduce an additional shot into the timeline. An automatically generated caption explicitly instructs the model to identify and remove this newly added shot, thus reverting the timeline to its original form. This direct instruction ensures the model learns the specific task through controlled adjustments. Each transformation is guided by instruction templates q_t , filled with the applicable cues c^i , to ensure the instructions are clear and relevant to the task at hand. We use this procedure to generate data with multiple lengths, and multiple instructions at a time that combine multiple atomic operations in one instruction, as we will show later in the experiment Section 4. We explain the data generation procedure in more detail in Section 3.2 of the supplementary material.

3.3 Training the Timeline Assembler

We illustrate our training procedure in Figure 3. The Timeline Assembler’s model $\mathcal{M}_{\gamma, \theta}$ has two learnable modules: the projection layer $h_\gamma(\cdot)$ and the language model $f_\theta(\cdot)$. Since we want to keep the multitask, instruction-following capabilities of the LLM, $f_\theta(\cdot)$, we keep its weights mostly frozen except for a lightweight set of learnable Low-Rank Adapters (LoRA) [15].

Our goal is to fine-tune the weights of the projection layer and large language model using our proxy assembly tasks gathered in the training dataset \mathcal{D} . In practice, we optimize the negative log-likelihood loss (NLL) as cited in Brown et al. (2020) by using the following: $\text{Loss}(\theta, \gamma) = -\sum_{i=1}^N \log P(\hat{\mathbf{X}}^i | \mathbf{X}^i; \theta, \gamma)$. where N is the number of samples and $P(\hat{\mathbf{X}}^i | \mathbf{X}^i; \theta, \gamma)$ is the probability assigned by the model $\mathcal{M}_{\gamma, \theta}$ to the correct output sequence $\hat{\mathbf{X}}^i$ given the input \mathbf{X}^i formed by the collection tokens \mathbf{X}_C^i , timeline tokens \mathbf{X}_S^i , and assembly instruction tokens \mathbf{X}_q^i for sample i .

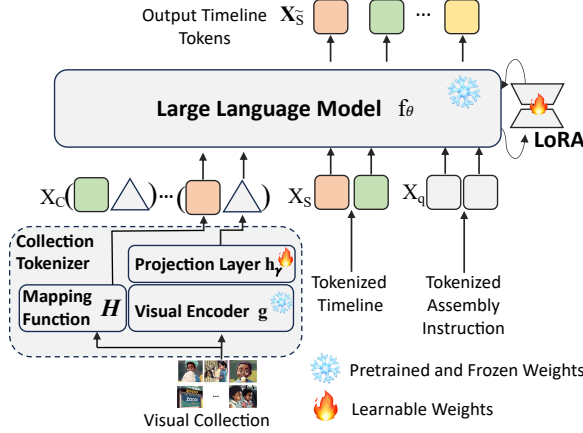


Figure 3: **Training the Timeline Assembler.** The Collection Tokenizer is composed of a frozen visual encoder $g(\cdot)$, a mapping function \mathcal{H} that generates an identifier token for each input visual asset (image/video) in the collection, and a learnable projection layer $h_\gamma(\cdot)$ that maps visual embeddings into visual tokens aligned with the LLM. We keep the LLM $f_\theta(\cdot)$ mostly frozen except for a lightweight set of learnable LoRA Adapters [15].

4 Experiments

Implementation Details. Our model has four main components as shown in Figure 3: the visual backbone, the projection layer, the LLM, and the low-rank matrices from LoRA [15]. As shown in Figure 3, the visual backbone and the LLM are kept frozen throughout the training. We only train the projection layer and the LoRA parameters. For the visual encoder, we adopt the same architecture as BLIP-2 [23] with ViT-g/14 from EVA-CLIP [11] and a Q-former that outputs 32 tokens per image [23]. We initialize the projection layer with that from MiniGPT-4 [54]. Finally, for the instruction-following LLM, we use Vicuna [53]. For the visual tasks, we use BLIP-2 Flan T5 XL [23] captions as cues c . We provide additional details in Section ?? in the supplementary material.

Evaluation Metrics. Our system’s performance is assessed using *assembly accuracy*, which determines whether a generated timeline exactly matches the corresponding ground-truth timeline. A true positive is a timeline with all its elements identical to the ground-truth, while a false positive occurs if any element in the predicted timeline differs from the ground-truth. We calculate the *Overall* assembly accuracy by evaluating each instance in the testing datasets and computing the percentage of correctly predicted samples. Additionally, we report the assembly accuracy for each cue type.

4.1 New Datasets for Instructed Visual Assembly

We present the details of our Instructed Assembly Datasets. For both cases, the size of the collection is 20 visual assets (images or videos). However, the design is not limited to this number (we show results with different collection sizes in the supplementary material). For the timeline length, we use $l = 5$ for Table 1, and $2 \leq l \leq 19$ for Section 4.4. Every task implies a single modification to the timeline. We use two data sources \tilde{D} to create Instruction Visual Assembly Datasets. Both are described below. To ensure instruction-phrasing generalization, we ensure that the training templates do not overlap with the validation and testing ones.

Visual Storytelling Assembly Dataset (VIST-A). We use images from [17] as source of image sequences to form our new Visual Storytelling Assembly dataset. The original dataset in [17], contains Flickr images linked to each other by annotators, to create visual stories each consisting of 5 images and their captions. We use the procedure explained in Section ?? to generate offline samples for each one of the assembly tasks. For testing, we create 80 samples per task for a total of 640 visual storytelling assembly tasks. For training, we create data online.

Video Sequence Assembly Dataset (VID-A). We collect a total of 12,088 YouTube Shorts [1]. We use these videos to show the capabilities of the Timeline Assembler for video assembly. We divide each video into shots and take the center frame as the shot representation. We create all possible timelines of a video by sampling every possible sequence of n consecutive shots. *The visual collection is formed by the shots of the video and when needed.* We construct the *Video Sequence Assembly Dataset (VID-A)* by creating 80 samples per task for a total of 640 video assembly tasks for testing. We create three validations sets using this procedure, one with timeline length of 5 for Section 4.2,

Table 1: **Instructed Visual Assembly Results.** We compare the Timeline Assembler against multiple baseline approaches. We include zero-shot results for powerful open-source and private VLMs. We compare the performance of our Assembler across various model capacities on two novel datasets, VST and VISTA, reporting assembly accuracy for Overall, Positional, and Semantic (cues). * denotes adjustment on the original implementation of the models.

	Assembly Accuracy(%)					
	VIST-A (Image-based)			VID-A (Video-based)		
	Positional	Semantic	Overall	Positional	Semantic	Overall
<i>Zero-shot</i>						
MiniGPT-4[54] *	0.0	0.0	0.0	0.0	0.0	0.0
LLaVA-1.5[8] *	5.6	0.0	2.8	4.7	0.3	2.5
GPT-4o	72.8	25.0	48.9	72.8	18.8	45.8
Ours						
Timeline Assembler-7B	90.1	58.1	74.1	91.8	41.8	66.8
Timeline Assembler-13B	96.9	66.4	81.6	93.7	47.5	70.6

one with variable timeline lengths for multi-len experiments in Section 4.4, and one with several instructions and variable timeline lengths for compositional experiments in Section 4.4. For training, we generate data samples using Algorithm ?? online.

4.2 Instructed Visual Assembly Results

In Table 1, we compare our model with state-of-the-art multimodal models. We observed that both MiniGPT-4 and LLaVA-1.5 lacked the capacity to manage a vast amount of visual data simultaneously. Accordingly, we adjusted these models to handle the entire collection at once; for more details on these adjustments, please refer to the supplementary material ???. After modifying both models for our task, neither showed satisfactory performance in these challenging tasks. We discuss the failure cases of these two models further in Section ?? of the supplementary material. We attribute their inadequate performance primarily to the data and instructions on which they were trained, which significantly differ from the task of instructed visual assembly.

To establish a stronger baseline, we evaluated the recent GPT-4o and found that this model achieved very impressive zero-shot performance on both datasets, recording 48.9% on VIST-A and 45.8% on VID-A. However, when trained specifically for the task, our Timeline Assembler clearly emerged as the best alternative for instructed visual assembly, outperforming all models in both positional and semantic tasks. The smaller Timeline Assembler-7B achieves 66.8% on the VID-A dataset. When scaled up to 13B parameters, the Timeline Assembler improves to 70.6%. Thus, our strategy proves effective and further enhances the capabilities of LLMs and multimodal LLMs for assembly tasks. For a detailed breakdown of the results for each task and additional analysis on the low performance of the baselines, please refer to Table ?? in the supplementary material. Additionally, we implemented additional baselines to perform instructed visual assembly using text-only language models, leveraging BLIP-2 [23] and replacing every image with its corresponding caption. Detailed explanations and results of these additional experiments are presented in Table ?? in the supplementary material.

4.3 Analysis of the Timeline Assembler

Impact of Multi-Task Training (Table 2). We compare the Timeline Assembler against specialized models trained for each one of the assembly tasks. We use a similar training strategy but train 8 individual models, one for each assembly task on the VIST-A dataset. To determine which of the 8 models to use for a given instruction, we propose three approaches:

- *Random Selection*: selects one of the eight models for a given instruction.
- *Oracle Task Classifier*: employs an Oracle Task classifier to choose the model.

Table 2: **Impact of Multi-Task Training.** We present a comparative analysis of our proposed model’s performance under single-task versus multi-task training paradigms. For the single-task case, we use one model per assembly task and a Task Classifier to decide which model to use. The Random and Oracle classifiers serve as lower-bound and upper-bound references for the single-task cases, respectively. We use GPT-4o to classify every instruction into one of the tasks. Our model (featured in the last row) is a single model that understands and performs every task and does not need any task classifier. For each task we report the assembly accuracy (%) on the VIST-A dataset.

	Multi Task Classifier	Task	Positional Cues				Semantic Cues				Avg.
			Insert	Remove	Replace	Swap	Insert	Remove	Replace	Swap	
Lower-Bound	✗	Random	10.0	15.0	17.5	11.3	20.0	17.5	15.0	16.3	15.3
Upper-Bound	✗	Oracle	99.2	98.8	100.0	98.8	71.3	70.0	69.6	26.3	79.2
Single-Task Models	✗	GPT-4o	88.1	97.6	80.0	76.5	67.7	70.0	46.1	12.1	67.3
Timeline Assembler	✓	N/A	98.8	90.4	85.4	85.8	65.8	66.7	54.2	45.8	74.1

- *GPT-4o Classification:* uses GPT-4o to classify the instruction and select the appropriate model. This classifier achieves an accuracy of 81.6%. Details about the performance of the GPT-4o Task classifier are provided in the supplementary material.

Table 2 shows the remarkable performance of the Timeline Assembler. It is surpassed by only 5% by the single-task models in conjunction with the Oracle classifier. This finding highlights the Timeline Assembler’s strong multi-tasking capabilities. The Timeline Assembler, being a single model that effectively handles multiple tasks without the need for a task classifier, proves to be highly practical. These results suggest that incorporating more assembly tasks could further enhance the Timeline Assembler model’s efficacy.

Impact of Training Scale. We analyze how data size impacts the Timeline Assembler’s learning. We limit the percentage of VIST-A and VID-A data samples available during training. Since we create training data on-the-fly, cutting down the number of sequences available during training means less variability in the data. Therefore, we scale the number of epochs accordingly. Table 3a shows three sizes for VIST-A and VID-A datasets. We observe that the availability of more data for creating assembly tasks on-the-fly leads to better performance.

Cross Dataset Analysis. Table 3b presents the results when our model is trained on one dataset and tested on the other. Our model demonstrates generalizability across datasets. When the Timeline Assembler is trained on visual stories, it can still learn how to perform video sequence assembly and vice versa. Complementary to the observations in Table 3a, when merging the two datasets during training, the Timeline Assembler performs the best across both datasets by gaining a notable 6% on VIST-A and 8% on VID-A. Thus, the favored approach is to continuously incorporate more data.

Ablation of Learnable Modules. Table 3c contrasts our model’s performance with and without its key components. First, when deactivating LoRA (first row), the Timeline Assembler’s performance drops by a significant 31%. Second, initializing the projection layer from [54] but freezing it (second row) instead of finetuning it, substantially lowers the assembly accuracy by 12.5%. Finally, when the projection layer is trained from scratch (third row), the performance is still competitive, as it only drops by 4%. These results validate that training both, the projection layer and LoRA adapters, is crucial in training Timeline Assembler to perform instructed visual assembly tasks.

4.4 Timeline Assembler Capabilities

In this section, we evaluate the Timeline Assembler on more complex tasks inspired by real video editing applications. We aim at evaluating the Timeline Assembler under two difficult scenarios: (i) The first scenario consists of handling variable input timeline lengths. This scenario challenges the Timeline Assembler with varying-length sequences and increasingly harder tasks as the length increases. (ii) The second scenario challenges the Timeline Assembler to deal with the composition of instructions specified in a single query. In this task, the model must understand multiple instructions at once and perform several timeline modifications to successfully resolve the user’s input.

Table 3: **Analysis of the Timeline Assembler.** (a) studies the impact of varying the training data size. (b) delves into the model’s generalization by testing its performance when trained on one dataset and tested on another; ‘mean’ is the average accuracy across both tests. (c) presents an ablation study, indicating the model’s performance with different variations of the model on VIST-A.

(a) Training scale			(b) Cross dataset analysis				(c) Learnable module ablation	
<i>Training scale</i>	VIST-A	VID-A	<i>Testing</i>				VIST-A	
10%	17.2	15.9	<i>Training</i>	VIST-A	VID-A	mean	w/o LoRA	43.0
50%	71.9	61.7	VIST-A	74.1	62.2	68.2	frozen $h_\gamma(\cdot)$ init. from [54]	61.6
100%	74.1	62.2	VID-A	71.7	66.8	69.2	$h_\gamma(\cdot)$ from scratch	70.5
			<i>all</i>	80.0	68.8	74.4	Ours	74.1

Multi-Length Timeline Assembler. In the previous sections, we evaluated the Timeline Assembler’s capabilities with fixed input timeline lengths of 5 assets. However, in real-life applications, dealing with multiple timeline lengths is essential. Therefore, using Algorithm ?? with the same data source as VID-A, we construct a test set for multiple-length assembly tasks, which we call multilen-VID-A. At training time, we also create data on-the-fly using Algorithm ?. Figure 4 reports the assembly accuracy of the assembler trained on multiple lengths against the strongest baseline GPT-4o from Table 1. We report the average assembly accuracy across all tasks within different ranges of timeline lengths (detailed per-task results can be found in the supplementary material). We observe a decrease in GPT-4o’s performance as the input timeline length increases to 5 and above, indicating that the tasks become more difficult. Our Timeline Assembler performs consistently across different lengths.

Compositional Timeline Assembler. Another desirable feature for the assembler is to handle compositions of assembly instructions. This compositional capability would allow greater flexibility in the assembler’s functionality. Therefore, we create a test set of compositional semantic assembly tasks, named Compositional-VID-A. Since semantic tasks are more challenging (as shown in Table 1), we combine them for our test set, presenting a highly challenging yet practical scenario for the assembler. Our test set includes combinations of two semantic tasks in a single instruction. Examples of such instructions can be found in Figures ?? to ?. In Table 4, we present different training strategies and their impact on assembly accuracy for multilen-VID-A and Compositional-VID-A. We note that GPT-4o struggles with multiple tasks simultaneously. Notably, the compositional Timeline Assembler outperforms GPT-4o and manages to perform compositional operations 36.3% of the time. Interestingly, the best way to train the Timeline Assembler for compositional tasks is to also incorporate single tasks (“Atomic-Task Training”) during training, as shown in rows 2 (compositional only) vs. 3 (compositional and atomic) of Table 4. The Compositional Timeline Assembler also performs well on Multi-len VID-A and VID-A, we report these results in the supplementary material.

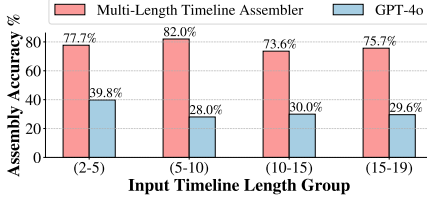


Figure 4: **Multi-Length Timeline Assembler.** Performance across various timeline lengths. Unlike GPT-4o, which has decreasing performance after timeline lengths larger than 5, our Timeline Assembler is consistent across different lengths.

5 Conclusion

We introduce the Timeline Assembler, the first generative model for instructed visual assembly. We train our model using automatically generated assembly tasks. We validate our approach on two newly built instruction visual assembly datasets and show that the Timeline Assembler follows assembly instructions more accurately than competing baselines including strong (multi-modal) large language models. Looking ahead, we believe the Timeline Assembler is a step towards generative models capable of complex reasoning (such as timeline editing or long-form story telling) over large collections of multi-modal assets. We include the limitations of our paper in ??.

Table 4: **Compositional Timeline Assembler.** We show results on compositional tasks. Our Compositional Timeline Assembler outperforms GPT-3.5. Additionally, we show the benefit of including single-tasks (atomic) during training.

Model	Training		Semantic Comp. VID-A (%)
	Atomic	Comp.	
GPT-4o	✗	✗	3.8
Timeline Assembler	✗	✓	26.3
Timeline Assembler	✓	✓	36.3

Acknowledgments The research reported in this publication was partially supported by funding from King Abdullah University of Science and Technology (KAUST) - Center of Excellence for Generative AI, under award number 5940.

References

- [1] YouTube Shorts - Creators. https://www.youtube.com/intl/en_in/creators/shorts/. Accessed: August 15, 2023. **6**
- [2] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. **3**
- [3] I. Arev, H. S. Park, Y. Sheikh, J. Hodgins, and A. Shamir. Automatic editing of footage from multiple social cameras. *ACM Transactions on Graphics (TOG)*, 33(4):1–11, 2014. **3**
- [4] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023. **3**
- [5] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. **3**
- [6] B. Chen, A. Ziai, R. S. Tucker, and Y. Xie. Match cutting: Finding cuts with smooth visual transitions. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2115–2125, 2023. **3**
- [7] G. Chen, Y.-D. Zheng, J. Wang, J. Xu, Y. Huang, J. Pan, Y. Wang, Y. Wang, Y. Qiao, T. Lu, et al. Videollm: Modeling video sequence with large language models. *arXiv preprint arXiv:2305.13292*, 2023. **3, 5**
- [8] W.-G. Chen, I. Spiridonova, J. Yang, J. Gao, and C. Li. Llava-interactive: An all-in-one demo for image chat, segmentation, generation and editing. 2023. **3, 5, 7**
- [9] F. Company. Google photos just got much better. here’s how [exclusive]. <https://www.fastcompany.com/90938324/google-photos-just-got-much-better-heres-how-exclusive>, 2024. Accessed: May 17, 2024. **2**
- [10] D. Driess, F. Xia, M. S. M. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, W. Huang, Y. Chebotar, P. Sermanet, D. Duckworth, S. Levine, V. Vanhoucke, K. Hausman, M. Toussaint, K. Greff, A. Zeng, I. Mordatch, and P. Florence. Palm-e: An embodied multimodal language model. In *arXiv preprint arXiv:2303.03378*, 2023. **3**
- [11] Y. Fang, W. Wang, B. Xie, Q. Sun, L. Wu, X. Wang, T. Huang, X. Wang, and Y. Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19358–19369, 2023. **6**
- [12] T.-J. Fu, X. E. Wang, S. T. Grafton, M. P. Eckstein, and W. Y. Wang. M3I: Language-based video editing via multi-modal multi-level transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10513–10522, 2022. **2**
- [13] S. Geng, J. Tan, S. Liu, Z. Fu, and Y. Zhang. Vip5: Towards multimodal foundation models for recommendation. *arXiv preprint arXiv:2305.14302*, 2023. **3**
- [14] Y. Hao, H. Song, L. Dong, S. Huang, Z. Chi, W. Wang, S. Ma, and F. Wei. Language models are general-purpose interfaces. *arXiv preprint arXiv:2206.06336*, 2022. **3**
- [15] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. **3, 5, 6**
- [16] S. Huang, Z. Jiang, H. Dong, Y. Qiao, P. Gao, and H. Li. Instruct2act: Mapping multi-modality instructions to robotic actions with large language model. *arXiv preprint arXiv:2305.11176*, 2023. **3**
- [17] T.-H. Huang, F. Ferraro, N. Mostafazadeh, I. Misra, A. Agrawal, J. Devlin, R. Girshick, X. He, P. Kohli, D. Batra, et al. Visual storytelling. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 1233–1239, 2016. **6**
- [18] B. Huber, H. V. Shin, B. Russell, O. Wang, and G. J. Mysore. B-Script: Transcript-based B-roll video editing with recommendations. In *ACM Conference on Human Factors in Computing Systems (CHI)*, 2019. **2**
- [19] Y. Jiang, A. Gupta, Z. Zhang, G. Wang, Y. Dou, Y. Chen, L. Fei-Fei, A. Anandkumar, Y. Zhu, and L. Fan. Vima: General robot manipulation with multimodal prompts. In *Fortieth International Conference on Machine Learning*, 2023. **3**
- [20] S. Kwon, J. Park, M. Kim, J. Cho, E. K. Ryu, and K. Lee. Image clustering conditioned on text criteria. *arXiv preprint arXiv:2310.18297*, 2023. **3**
- [21] M. Leake, A. Davis, A. Truong, and M. Agrawala. Computational video editing for dialogue-driven scenes. *ACM Trans. Graph.*, 36(4):130–1, 2017. **3**
- [22] B. Li, Y. Zhang, L. Chen, J. Wang, J. Yang, and Z. Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023. **3**
- [23] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. **6, 7**
- [24] K. Li, Y. He, Y. Wang, Y. Li, W. Wang, P. Luo, Y. Wang, L. Wang, and Y. Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. **3**

- [25] Y. Z. B. G. N. Li and Y. Z. Q. W. Z. Yu. Representation learning of next shot selection for vlog editing. *CVEU Workshop at ICCV*, 2023. 3
- [26] K. Lin, C. Agia, T. Migimatsu, M. Pavone, and J. Bohg. Text2motion: From natural language instructions to feasible plans. *arXiv preprint arXiv:2303.12153*, 2023. 3
- [27] K. Lin, F. Ahmed, L. Li, C.-C. Lin, E. Azarnasab, Z. Yang, J. Wang, L. Liang, Z. Liu, Y. Lu, et al. Mm-vid: Advancing video understanding with gpt-4v (ision). *arXiv preprint arXiv:2310.19773*, 2023. 5
- [28] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 2, 3, 5
- [29] R. Liu, C. Li, Y. Ge, Y. Shan, T. H. Li, and G. Li. One for all: Video conversation is feasible without video instruction tuning. *arXiv preprint arXiv:2309.15785*, 2023. 3
- [30] Y. Lu, F. Ni, H. Wang, X. Guo, L. Zhu, Z. Yang, R. Song, L. Cheng, and Y. Yang. Show me a video: A large-scale narrated video dataset for coherent story illustration. *IEEE Transactions on Multimedia*, 2023. 3
- [31] M. Maaz, H. Rasheed, S. Khan, and F. S. Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023. 3
- [32] F. Manjoo. Google photos is your perfect jukebox for all your memories. *The New York Times*, 11 2018. Accessed: May 17, 2024. 2
- [33] J. Merullo, L. Castricato, C. Eickhoff, and E. Pavlick. Linearly mapping from image to text space. *arXiv preprint arXiv:2209.15162*, 2022. 4
- [34] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022. 3, 5
- [35] A. Pardo, F. Caba, J. L. Alcázar, A. K. Thabet, and B. Ghanem. Learning to cut by watching movies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6858–6868, 2021. 2, 3
- [36] S. Pei, J. Yu, Q. Chen, and W. He. Automatch: A large-scale audio beat matching benchmark for boosting deep learning assistant video editing. *arXiv preprint arXiv:2303.01884*, 2023. 2, 3
- [37] B. Qin, J. Li, S. Tang, T.-S. Chua, and Y. Zhuang. Instructvid2vid: Controllable video editing with natural language instructions. *arXiv preprint arXiv:2305.12328*, 2023. 2
- [38] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 3
- [39] A. Rao, X. Jiang, S. Wang, Y. Guo, Z. Liu, B. Dai, L. Pang, X. Wu, D. Lin, and L. Jin. Temporal and contextual transformer for multi-camera editing of tv shows. *arXiv preprint arXiv:2210.08737*, 2022. 3
- [40] A. Rao, J. Wang, L. Xu, X. Jiang, Q. Huang, B. Zhou, and D. Lin. A unified framework for shot type classification based on subject centric lens. In *The European Conference on Computer Vision (ECCV)*, 2020. 3
- [41] R. Sennrich, B. Haddow, and A. Birch. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*, 2015. 4
- [42] Y. Shen, L. Zhang, K. Xu, and X. Jin. Autotransition: Learning to recommend video transition effects. In *European Conference on Computer Vision*, pages 285–300. Springer, 2022. 3
- [43] E. Song, W. Chai, G. Wang, Y. Zhang, H. Zhou, F. Wu, X. Guo, T. Ye, Y. Lu, J.-N. Hwang, et al. Moviechat: From dense token to sparse memory for long video understanding. *arXiv preprint arXiv:2307.16449*, 2023. 3
- [44] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 5
- [45] A. Truong, F. Berthouzoz, W. Li, and M. Agrawala. Quickcut: An interactive tool for editing narrated video. In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology*, pages 497–507, 2016. 2
- [46] M. Wang, G.-W. Yang, S.-M. Hu, S.-T. Yau, A. Shamir, et al. Write-a-video: computational video montage from themed text. *ACM Trans. Graph.*, 38(6):177–1, 2019. 3
- [47] Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, D. Khashabi, and H. Hajishirzi. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*, 2022. 2, 3
- [48] Z. Wang, J. Li, and Y.-G. Jiang. Story-driven video editing. *IEEE Transactions on Multimedia*, 23:4027–4036, 2020. 3
- [49] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021. 2, 3
- [50] Y. Xiong, F. C. Heilbron, and D. Lin. Transcript to video: Efficient clip sequencing from texts. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 5407–5416, 2022. 2, 3
- [51] Z. Yin, J. Wang, J. Cao, Z. Shi, D. Liu, M. Li, L. Sheng, L. Bai, X. Huang, Z. Wang, et al. Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. *arXiv preprint arXiv:2306.06687*, 2023. 3
- [52] H. Zhang, X. Li, and L. Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 3

- [53] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. P. Xing, H. Zhang, J. E. Gonzalez, and I. Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. [6](#)
- [54] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. [2](#), [3](#), [6](#), [7](#), [8](#), [9](#)