



# LLAVIDAL : Benchmarking Large LAnguage VIision Models for Daily Activities of Living

Rajatshubhra Chakraborty\* Arkaprava Sinha\* Dominick Reilly\* Manish Kumar Govind

Pu Wang Francois Bremond† Srijan Das

UNC Charlotte †Inria †Université Côte d’Azur

\* Equal contribution {rchakra6, asinha13, dreilly1}@charlotte.edu

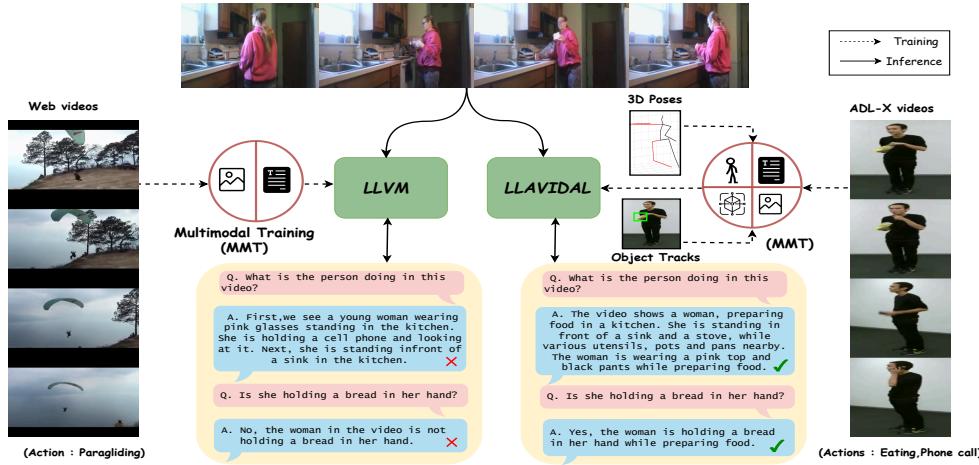


Figure 1: **Comparison of LLVM vs LLAVIDAL** : In real world scenarios, web-video trained models struggle to understand Activities of Daily Living due to the subtle nuances in the video, whereas our ADL-X trained LLAVIDAL model triumphs in understanding complex human-object interactions.

## Abstract

Large Language Vision Models (LLVMs) have demonstrated effectiveness in processing internet videos, yet they struggle with the visually perplexing dynamics present in Activities of Daily Living (ADL) due to limited pertinent datasets and models tailored to relevant cues. To this end, we propose a framework for curating ADL multiview datasets to fine-tune LLVMs, resulting in the creation of **ADL-X**, comprising 100K RGB video-instruction pairs, language descriptions, 3D skeletons, and action-conditioned object trajectories. We introduce **LLAVIDAL**, an LLVM capable of incorporating 3D poses and relevant object trajectories to understand the intricate spatiotemporal relationships within ADLs. Furthermore, we present a novel benchmark, **ADLMCQ**, for quantifying LLVM effectiveness in ADL scenarios. When trained on ADL-X, LLAVIDAL consistently achieves state-of-the-art performance across all ADL evaluation metrics. Qualitative analysis reveals LLAVIDAL’s temporal reasoning capabilities in understanding ADL. The link to the dataset is provided at: <https://adl-x.github.io>/

## 1 Introduction

Human cognitive perception integrates information from multiple sensory modalities to form a unified representation of the world [1]. Towards emulating human cognitive perception in digital intelligence,

initial efforts focused on integrating vision and language modalities [2, 3, 4, 5, 6]. Subsequently, the success of LLMs like GPT [7], PALM [8], BLOOM [9] led to the introduction of multimodal conversational models[10, 11, 12, 13, 14, 15, 16] that combine image pixels and LLMs, we dub as Large Language-Vision Language Models (LLVMs). However, these image-based LLVMs lack the capability for complex reasoning and interactions, particularly in understanding spatio-temporal relationships involved in human activities. In this study, we investigate the understanding of Activities of Daily Living (ADL) videos by LLVMs, which present various challenges including multiple exocentric viewpoints, fine-grained activities with subtle motion, complex human-object interactions, and long-term temporal relationships. We envision that LLVMs capable of addressing these challenges will significantly influence the future intelligent systems, particularly in healthcare applications such as eldercare monitoring, cognitive decline assessment, and robotic assistance development.

Recently, [17, 18, 19, 20, 21, 22, 23] have integrated videos into LLMs, leading to the development of video-based LLVMs capable of capturing spatio-temporal features. However, these models are predominantly trained on large-scale web videos [24, 25, 26, 27, 28], which mainly consists of sports clips, movie excerpts, and instructional videos. These videos, typically filmed by professionals, follow strict temporal sequences in closely controlled background (e.g., Paragliding). The evident temporal structure and scene semantics in such videos facilitate spatial understanding within LLVMs, as shown in 1. In contrast, ADL videos pose additional challenges, characterized by temporal unstructuredness where diverse actions may unfold concurrently within a single sequence [29]. For instance, *a person cooking could intermittently engage in unrelated activities like making a phone call or drinking water, disrupting the linear progression of the composite action cooking*. Consequently, existing LLVMs trained on web videos struggle to capture such visually perplexing dynamics inherent in ADL scenarios. Moreover, unlike specialized video architectures designed for understanding ADL [30, 31, 32, 33, 34, 35, 36], these LLVMs lack explicit utilization of cues like 3D poses or object encodings, which are crucial for understanding ADL. These cues aid in learning view-invariant representations and capturing fine-grained details essential for interpreting complex human activities. Hence, the current limitations in understanding ADL stem from the lack of instruction tuning of LLVMs on real-world multiview ADL datasets captured in indoor settings and the simplistic design of LLVMs with holistic operations.

To this end, we propose a framework of curating ADL videos for instruction tuning LLVMs. This framework introduces the **ADL-X** dataset, comprising 100K untrimmed RGB video-instruction pairs, 3D poses (P), language descriptions, and action-conditioned object trajectories (see Table 1). We then introduce the **Large LAnguage VIision model for Daily Activities of Living (LLAVIDAL)**, trained on ADL-X, which integrates videos, 3D poses, and object cues into the LLM embedding space. Our study explores various strategies for integrating 3D pose information and human-object interactions within LLVMs, demonstrating that language contextualized features extracted from 3D poses and object trajectories can effectively be integrated into LLAVIDAL. Furthermore, we introduce a benchmark ADL Multiple Choices Question (**ADLMCQ**), specifically designed to evaluate the effectiveness of LLVMs for ADL. ADLMCQ includes action recognition (ADLMCQ-AR) and action forecasting (ADLMCQ-AF), assessed through a multiple choice question-answering task. We also evaluate existing LLVMs for generating video description of ADL scenes and compare their performance with LLAVIDAL. Our empirical findings indicate that LLAVIDAL with object cues, outperforms other LLVMs, including those trained on datasets of ten times the size, on the ADL benchmarks.

To summarize our contributions:

- We introduce ADL-X, the first multiview RGBD instruction ADL dataset, curated through a novel semi-automated framework for training LLVMs.
- LLAVIDAL is introduced as the first LLVM tailored for ADL, incorporating 3D poses and object cues into the embedding space of the LLM.
- A new benchmark, ADLMCQ, is proposed for an objective evaluation of LLVMs on ADL tasks, featuring MCQ tasks for action recognition & forecasting.
- Exhaustive experiments are conducted to determine the optimal strategy for integrating poses or objects into LLAVIDAL. Evaluation of existing LLVMs on ADLMCQ and video description tasks reveals that LLAVIDAL trained on ADL-X significantly outperforms baseline LLVMs.

Table 1: Video Instruction Dataset Comparison.

Dataset	Modalities	Subjects	Multiple Views	Videos	QA Pairs	Atomic Actions per Vid	Temporal Rand.	Object Traj.	Type
TimeIT[21]	RGB+L	NA	No	173000	173K	Medium	No	No	Web
VideoChat[17]	RGB+L	NA	No	8196	11K	Low	No	No	Web
Valley[26]	RGB+L	NA	No	64,687	65K	Low	No	No	Web
VideoChatGPT [20]	RGB+L	NA	No	27,801	100K	Medium	No	No	Web
<b>ADL-X</b>	<b>RGB+P+L</b>	<b>106</b>	<b>Yes</b>	<b>16,343</b>	<b>100K</b>	<b>High</b>	<b>Yes</b>	<b>Yes</b>	<b>ADL</b>

## 2 Semi-automated Framework for generating ADL Video-instructions Pairs

This section describes the data curation framework employed for the creation of a novel dataset, ADL-X. This dataset specifically caters to the instruction tuning of LLVMs within the ADL domain. ADL-X comprises video recordings of ADLs. To enrich the dataset and facilitate LLM training, question-answer (QA) pairs were generated from a corpus of long-form ADL videos. These QA pairs target various aspects of the ADLs, including: human pose configuration, objects relevant to the human actions, scene appearance, and the fine-grained actions performed. We hypothesize that incorporating such instructional tuning during the LLVM training process will promote alignment of visual tokens within the LLM’s embedding space. ADL-X represents a comprehensive ADL dataset encompassing various modalities: - RGB videos, 3D poses, Language descriptions, object tracklets. This rich dataset offers a valuable tool for evaluating the capabilities of LLVMs in tasks related to ADLs, including description, recognition, and anticipation.

A critical characteristic of ADL videos lies in the inherent spontaneity of the actions performed. Unlike scripted scenarios [25, 37, 38], fine-grained actions within ADLs often occur randomly. To capture this essential characteristic within our dataset, we curated ADL-X from NTU RGB+D 120 dataset [39]. This selection was motivated by the dataset’s focus on ADL videos and its inherent diversity in terms of actions, subjects, and camera viewpoints. Also, this data curation framework could be extended to any existing trimmed/untrimmed ADL datasets [40, 41, 42]. Below, we elaborate the steps involved in building the ADL-X in a chronological order.

**Person-centric Cropping.** ADL tasks necessitate a focus on the individual performing the actions, the actions themselves, and the human-object interactions. To achieve this targeted focus within the data curation framework, we implemented a person-centric cropping strategy leveraging the pose information captured through Kinect sensors [43]. By using the pose information in each frame of the NTU RGB+D 120 dataset, we are able to detect and crop out the person(s) performing the actions. This cropping process effectively reduces the amount of background information present in the videos, eliminating data irrelevant to the target ADLs. This step is crucial as existing ADL datasets often contain extensive background information that is not relevant to the actions being performed. The presence of such extraneous information can significantly hinder subsequent stages within the data curation framework.

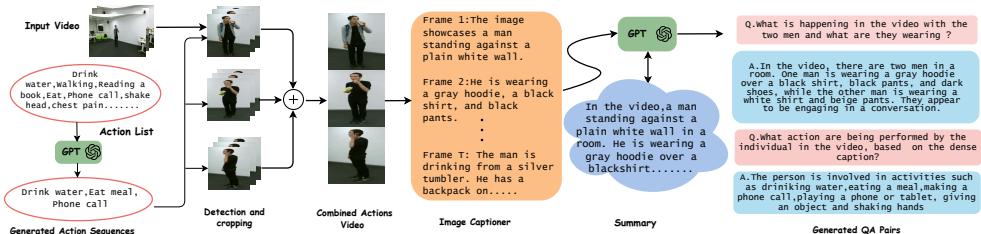


Figure 2: Dataset Curation Pipeline: We employ CogVLM[44] as our person-centric image captioner and GPT 3.5 Turbo[7] as our summarizer and QA generator.

**Stitching shorts clips.** To capture the inherent randomness of real-world ADLs, we constructed a set of 160 composite action sequences. These sequences were generated by prompting a GPT to combine individual actions from the original NTU RGB+D 120 dataset’s list of 120 actions (denoted as  $A_1, A_2, \dots, A_{120}$ ). An example sequence structure could be represented as  $A_1 \rightarrow A_3 \rightarrow A_{17}$ . Following these generated composite action sequences, we temporally stitched together short video clips ( $clip_j^a$ , where  $a$  is the action class) from the NTU dataset. This stitching process ensured that all clips within a video belonged to the same subject and camera view, maintaining coherence in the resulting video sequence. For instance, a stitched video sequence might be represented as  $[clip_{r1}^1 \ clip_{r2}^3 \ clip_{r3}^{17}]$  where  $r1, r2,$

$r_3$  represent unique clip identifiers within the dataset for the specific subject performing the actions (actions 1, 3, and 17, respectively). The intentional randomness of the generated action sequences reflects the unstructured flow of actions encountered in ADL. To further enhance diversity and ensure no bias towards specific subject-action combinations, we shuffled both the action sequences and the subject assignments. This process resulted in the creation of **16,343 stitched videos** with an average 5 actions per video.

**Frame Level Captioning and Dense Descriptions.** This step is the process of generating weak pseudo-labels for automated instruction tuning of the LLVM with the curated dataset. An image captioning model CogVLM [44] is employed to automatically generate frame-level captions for the stitched ADL videos at a rate of  $0.5\text{fps}$ . These captions are subsequently compiled into a dictionary linking each frame identifier to its corresponding description. To enhance the reliability of the pseudo-labels, we implemented an action-conditioned filtering while generating the video descriptions. The dictionary with the frame descriptions, along with the action labels present in the stitched videos, are then used to prompt a GPT 3.5 turbo model to generate a cohesive structured description of the entire stitched video, constrained to a maximum of 300 words. This step leverages the known action labels associated with each video to remove irrelevant noise potentially introduced during the caption generation process. We evaluated various image captioning models, including BLIP-2 [45], and InstructBLIP [46] for frame-level caption generation. However, CogVLM is ultimately chosen due to its ability to generate denser and appropriate descriptions. Please refer to Appendix H for our detailed prompting strategy in generating the descriptions.

**Generating QA Pairs.** LLVMs necessitate training data in the form of question-answer (QA) pairs. To generate domain-specific QA pairs for ADL, we leverage the dense video descriptions obtained in the previous step as illustrated in Figure 2. An instruction template (detailed in Appendix H) guides GPT-3.5 in formulating questions across various categories relevant to ADL. These categories include: video summary, performed actions, spatial details, human-object interactions and other video-specific inquiries. Through this prompting approach, we curate a dataset of **100K video instruction pairs**, namely ADL-X, for the stitched ADL videos. These QA pairs benefit from the detailed descriptions and person-centric cropping, resulting in reduced LLM hallucinations compared to other existing methods [17, 20].

Notably, the framework employed for constructing ADL-X from trimmed, labeled action videos can be generalized to other existing datasets. This generalization paves the way for efficient training of domain-specific LLVMs.

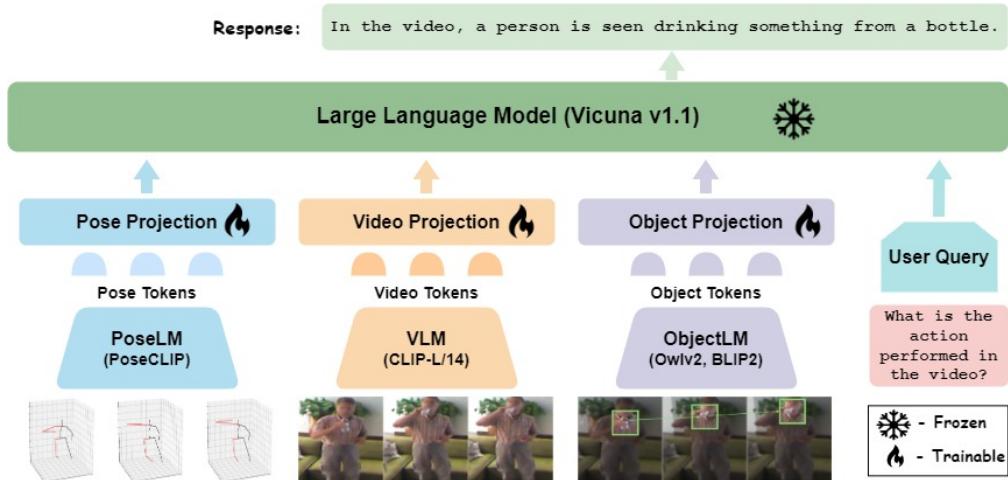


Figure 3: Overview of **LLAVIDAL**, which utilizes an LLM to integrate multiple modalities, including video, pose, and object features. Videos are represented by embeddings obtained from a **VLM**, poses are processed through (**PoseLM**), and object embeddings are obtained through (**ObjectLM**). These embeddings are projected into the LLM space, where they are concatenated with tokenized text queries for instruction tuning.

### 3 LLAVIDAL: An LLVM for ADL

LLAVIDAL is a large language vision model designed to align ADL videos with an LLM to generate meaningful conversation about the daily activities performed by humans. This model, similar to Video-ChatGPT [20] and LLaVA [18], integrates a visual encoder with the Vicuna language decoder [47] and is fine-tuned on instructional language-vision data. Unlike Video-ChatGPT [20] and LLaVA [18], LLAVIDAL leverages the random temporal structure present in ADL-X and incorporates additional data modalities such as 3D human poses and human-object interaction cues. This allows LLAVIDAL to generate accurate conversations that are not only contextually appropriate but also temporally aligned with the human activities depicted in the input video. This section will first present a background of LLVM models to align videos with LLMs. Then, we will outline the strategies employed to integrate 3D poses and object interaction cues within the language space of the LLM for enhanced understanding of videos featuring ADL. Subsequently, we will describe the training architecture of LLAVIDAL.

#### 3.1 Background: LLVM

Following [20], given an input video denoted by  $\nu_i \in \mathbb{R}^{T \times H \times W \times C}$ , where  $T$  represents the frames encoded using a pretrained vision-language model (**VLM**) CLIP-L/14 [2] to obtain frame-level embeddings for the video,  $x_i \in \mathbb{R}^{T \times h \times w \times D}$ , with  $D$  as the embedding dimension, and  $h = H/p$ ,  $w = W/p$  representing the dimensions adjusted by patch size  $p$ . Temporal and spatial features are extracted by aggregating these frame-level embeddings along the respective dimensions. The video-level features,  $V_i \in \mathbb{R}^{F_v \times D_v}$ , are obtained by concatenating the temporal and spatial features, where  $F_v$  represents the spatio-temporal tokens and  $D_v$  is the video feature dimension. The video features are projected into the LLM embedding space using a linear projection layer  $\mathcal{T}_v$ . Thus, we obtain input tokens  $Q_v$  for the video features:

$$Q_v = \mathcal{T}_v(V_i) \in \mathbb{R}^{F_v \times K} \quad (1)$$

The text query is also tokenized such that  $Q_t \in \mathbb{R}^{F_t \times K}$ . The text query  $Q_t$ , refers to a question from the training data. The input to the LLM is the concatenation of  $Q_t$  and  $Q_v$  following the template : [USER:  $\langle Q_t \rangle \langle Q_v \rangle$  Assistant: ]. We perform instruction-tuning of the LLM on the prediction tokens, using its original auto-regressive training objective. The parameters of the LLM are frozen, thus the loss gradients only propagate through the projection layer  $\mathcal{T}_v$ .

#### 3.2 3D Poses for LLAVIDAL

ADL are rich in actions that primarily involve the movements of critical body parts or joints. The dataset ADL-X includes 3D human poses, which can be utilized to incorporate human kinematics and view-invariant features into the input embedding space of a LLM. These poses can be integrated into the LLM input space in several ways: as an additional text query  $Q_t$  for instruction tuning of the LLM, by deriving language descriptions of joint movements to provide context for the LLM, or through features extracted using a suitable pose-language encoder.

**Poses as QA.** We input the 3D joint coordinates alongside the associated human action from the video into GPT-3.5 Turbo [7], which generates a general description of the pose. This description is then re-fed into GPT-3.5 Turbo to generate two QA pairs that provide detailed explanations of the action’s motions. These QA pairs are subsequently added to the set of text queries  $Q_t$  in our training set for instruction tuning the LLM.

**Poses as Context.** To extract contextual information from human poses, we initially identify five peripheral joints — the head, right hand, left hand, right knee, and left knee — due to their significant contribution to motion in various actions. Using GPT-3.5 Turbo, we generate descriptions of the motion for each of these joints based on their trajectories throughout the video, specifically focusing on how the coordinates of these five joints evolve. The generated descriptions, denoted as  $Q_t^p$ , are subsequently appended to the text query  $Q_t$ , incorporates these pose descriptions as additional contextual information. This enriched query  $Q_t^{new} = [Q_t^p \ Q_t]$  is then employed for instruction tuning of the LLAVIDAL.

**Poses as Features.** To incorporate poses as tokens into the LLM, it is crucial to align the pose features with a language-contextualized space. To achieve this, we utilize a pretrained Pose-Language model (**PoseLM**), specifically PoseCLIP, to extract pose features that are aligned with the language domain. The PoseCLIP model comprises a pose backbone [48] and a CLIP text encoder [2], and it undergoes training in two phases. Initially, the pose backbone is pretrained on the NTU RGB+D dataset [49] for action classification. Subsequently, in the second phase, we optimize the similarity between pose features and text features, which encode the prompts describing their action labels,

using cross-entropy supervision as outlined in [3]. Further details on the training of this model are provided in Appendix C. These pose features, denoted as  $P_i \in \mathbb{R}^{F_p \times D_p}$ , where  $D_p$  represents the pose feature dimension, can be utilized as input tokens for training LLAVIDAL.

### 3.3 Action-Conditioned Object Cue for LLAVIDAL

To comprehensively understand ADL, it is crucial to not only grasp the semantics of objects but also their trajectories, which are closely linked to the actions performed. Consequently, we propose to explicitly utilize these object trajectories as integral components for training LLAVIDAL. Our framework involves a two-stage pipeline to extract object information directly from RGB video data: (i) *Action-conditioned object detection* and (ii) *Object Localization and Tracking*. Both stages leverage off-the-shelf models that are effective without the need for additional training, facilitating integration into LLAVIDAL for ADL analysis.

**Action conditioned object detection.** Given a stitched ADL video, which comprises a sequence of trimmed video segments (denoted as  $clip_j$ ), the first stage extracts the categories of objects present that are pertinent to the actions performed within each clip. We uniformly sample 8 frames from each video and employ a pre-trained BLIP-2 model [45] to generate a list of distinct objects observed in the frames. To avoid training LLAVIDAL with noisy data, we perform a filtering on the list of objects using the ground-truth action labels and GPT-3.5. Specifically, for each  $clip_j$  within a stitched video, we input the corresponding action label and the list of detected objects to GPT-3.5 and prompt it to identify the object(s) most relevant to the given action. For instance, if the objects *plant*, *chair*, *bottle*, *table* are detected in a video labeled with the action *Drinking*, GPT-3.5 is expected to filter out and select [*bottle*] as the relevant object for  $clip_j$ . Refer to Appendix H for our detailed action conditioned object detection prompting strategy.

**Object Localization and Tracking.** Given the list of relevant objects identified in the first stage, the second stage involves spatial localization of these objects within the scene and their temporal association (i.e., object tracking) based on the feature similarity of the image regions corresponding to the localized objects in the stitched video. We employ a pre-trained open vocabulary object localization model (**ObjectLM**), OWLv2 [50], and input the list of relevant objects detected in stage 1 along with the corresponding video. Localization and tracking are performed on 8 frames that are uniformly sampled from  $clip_j$  within a stitched video. For each frame, we obtain bounding boxes  $B_t \in \mathbb{R}^{n \times 4}$ , where each bounding box corresponds to one of the  $n$  relevant objects in the  $t$ th frame. Features for each object are then extracted from the image regions within these bounding boxes using our object localization model. We denote the features for the objects in frame  $t$  as  $O_t \in \mathbb{R}^{8n \times D_o}$ , where  $D_o$  is the object feature dimension. To associate objects across frames, we utilize a feature-based object tracking approach. Specifically, for each object in frame  $t$ , represented by the feature vector  $O_i^t \in \mathbb{R}^{D_o}$ , we compute the cosine similarity between  $O_i^t$  and all feature vectors in frame  $t + 1$ . The object  $i$  in frame  $t$  is then associated with the object in frame  $t + 1$  that exhibits the highest similarity score. This matching process is iterated for all objects in each frame, thereby establishing a track for each relevant object throughout the sampled frames. These object tracks, with corresponding bounding boxes and features, facilitate the integration of object information into the training of LLAVIDAL: Object as QA, Object as context, and Object as features.

**Object as QA.** Similar to the approach taken with poses, to generate QA pairs for objects, we formulate a question based on the trajectory coordinates of the relevant object(s). These QA pairs are added to the set of text queries  $Q_t$  for instruction tuning LLAVIDAL.

**Object as Context.** To integrate the context of detected objects into the LLM space, we append the list of relevant object labels, denoted by  $Q_t^o$ , to each text query token  $Q_t$ . Consequently, the updated text query is represented as  $Q_t^{new} = [Q_t^o \ Q_t]$ . This enhanced text query,  $Q_t^{new}$ , is utilized for instruction tuning.

**Object as Features.** The object features extracted during the object localization and tracking stage are utilized as input tokens  $Q_o \in \mathbb{R}^{8n \times D_o}$ , which are incorporated alongside the text query tokens ( $Q_t$ ) and input video tokens ( $Q_v$ ). For  $n$  relevant objects detected, the object query  $Q_o$  is structured using the following template  $[\langle Q_o \rangle = \langle Q_o^1 \rangle \langle Q_o^2 \rangle \dots \langle Q_o^n \rangle]$  where  $Q_o^j \in \mathbb{R}^{8 \times D_o}$  represent the features of each relevant object in the video.

### 3.4 Training LLAVIDAL

As illustrated in Figure 3, the QA pairs, along with context or features obtained from the RGB video, 3D poses, and object cues can be integrated into LLAVIDAL. Integrating QA pairs and contextual information is straightforward; they are introduced into  $Q_t$  and trained using standard methods for LLVM. However, to integrate other modalities with features, we feed these additional cues through specific projection layers designed to align them with the input space of the LLM. Accordingly, the

Table 2: Impact of ADL-X Training

Method	Training Data	ADLMCQ-AR (Smarthome)	ADLMCQ-AF (LEMMA)	Action Object	Description Action	(Charades) Correctness
VideoChatGPT [20]	ActivityNet	40.8	35.7	14.8	<b>16.1</b>	35.8
VideoChatGPT [20]	NTU120	49.8	33.5	27.0	10.1	38.8
ADL-X ChatGPT [20]	ADL-X	<b>52.3</b>	<b>44.8</b>	<b>32.2</b>	13.4	<b>43.0</b>

video, pose, and object features are projected into the LLM embedding space using linear projection layers  $\mathcal{T}_j$  for each cue  $j = \{v, p, o\}$ , resulting in LLM input token representation of the video, pose, and object cues, respectively:

$$Q_v = \mathcal{T}_v(V_i); \quad Q_p = \mathcal{T}_p(P_i); \quad Q_o = \mathcal{T}_o(O_i) \quad (2)$$

where  $Q_j \in \mathbb{R}^{F_j \times K}$ . Thus, the input to the LLM comprises the concatenation of  $Q_t$  and  $Q_j$  for  $j = \{v, p, o\}$ , structured according to the template: [USER:  $\langle Q_t \rangle \langle Q_v \rangle \langle Q_o \rangle \langle Q_p \rangle$  Assistant: ]. This training scheme ensures that the video, object, and pose cues are effectively aligned to the LLM embedding space, facilitating an accurate understanding of ADL. During the **inference**, LLAVIDAL utilizes only the holistic video cue, omitting person-centric cropping and consequently eliminating additional cues. In practice, the embedding dimensions are  $D_v = 1024$  for visual,  $D_o = 512$  for object features,  $D_p = 216$  for pose features and  $K = 4096$ . The number of tokens is set as  $F_v = 356$  and  $F_p = 256$  for visual and pose tokens respectively. We train LLAVIDAL for 3 epochs with a batch size of 32 and a learning rate of  $2e^{-5}$  on 8 A6000 48GB GPUs. For the purpose of promoting research in this field, we also provide the pose features and object trajectories of LLAVIDAL along with the dataset.

## 4 Experiments

### 4.1 Experimental Setting

**Evaluation Metrics.** Inspired by [20], LLVM’s ability to generate video-level descriptions is evaluated. This involves comparing the generated descriptions with ground truth and scoring them on dimensions such as Correctness of Information, Detail Orientation, Contextual Understanding, Temporal Understanding, and Consistency, with scores scaled to be bounded at 100. Due to the subjective nature of this metric, Mementos Evaluation [51] is also conducted to assess the recognition of common action-verbs and object-nouns in the video descriptions compared to ground truth, presenting F1 scores for these classifications. However, comparing video descriptions generated by LLVMs presents a challenge due to the inherently subjective nature of these descriptions. Some objective evaluation benchmarks for LLVMs [52, 53, 54] primarily focus on video tasks involving in-the-wild activities. Therefore, this paper introduces novel benchmarks for assessing LLVM’s temporal understanding of ADL videos. We propose two new **ADLMCQ** benchmarks including ADLMCQ-AR and ADLMCQ-AF. ADLMCQ-AR involves multiple-choice question-answering for action recognition, where the model selects the correct action from a set of options given a question about the action performed in a video. Similarly, ADLMCQ-AF focuses on action forecasting, requiring the model to predict the next action based on the preceding actions. It is important to note that all evaluations are performed zero-shot.

**Evaluation Datasets.** For ADLMCQ-AR evaluation, we utilize the Charades [55] and Toyota Smarthome [56] datasets. Evaluation for ADLMCQ-AF is conducted using LEMMA [57] and Toyota Smarthome Untrimmed (TSU) [58] datasets. Video description tasks are assessed using the Charades and TSU datasets, both featuring long-duration videos with multiple actions per video. Notably, for the TSU dataset, we manually annotated video descriptions with fine-grained details regarding activities performed by elderly individuals, employing 6 human annotators for 174 videos. Our evaluation relies on these annotated descriptions, which we also provide to the community as part of the test set for ADL-X.

### 4.2 Impact of ADL-X Training on LLVMs

To understand the requirement of ADL-X, we assess VideoChatGPT [20] trained on 100K instruction pairs from ActivityNet [25], trimmed NTU120 [39], and ADL-X in Table 2. Notably, ADL-X ChatGPT, trained on ADL-X, consistently outperforms the others in both ADLMCQ-AR and ADLMCQ-AF tasks. However, it’s worth mentioning that while the baseline [20] exhibits strong performance in the action metric of Mementos, it notably underperforms in the object metric. It’s important to emphasize that ADLMCQ evaluations offer more objective and reliable assessments for understanding the temporal comprehension of LLVMs.

Table 3: Introducing Pose and Object Cues into LLAVIDAL

Method	ADLMCQ-AR		ADLMCQ-AF		AD (Charades)		AD (TSU)	
	Charades	Smarthouse	LEMMA	TSU	Object	Action	Object	Action
ADL-X ChatGPT	58.0	52.3	44.8	25.25	16.6	14.8	16.6	14.8
Pose QA	48.5	49.0	42.0	21.2	31.8	14.0	16.5	15.9
Pose Context (PC)	50.8	54.0	45.0	22.3	30.5	<b>14.8</b>	18.6	15.4
Pose Features (PF)	56.7	<b>57.0</b>	<b>51.3</b>	26.0	<b>32.7</b>	13.5	18.2	13.0
PC + PF	52.5	53.1	44.6	24.9	32.1	13.6	17.5	15.6
Object QA	51.1	50.1	40.3	23.0	32.1	13.7	17.0	16.0
Object Context	44.6	46.2	41.8	21.0	31.2	<b>14.7</b>	17.2	16.5
Object Features (OF)	<b>59.0</b>	<b>58.8</b>	<b>52.6</b>	<b>27.0</b>	<b>33.1</b>	14.3	18.0	<b>17.7</b>
PF + OF	56.2	56.1	51.0	26.6	30.4	14.1	<b>20.0</b>	14.1

### 4.3 How to introduce object and pose cues into the LLM space?

Table 3 explores the integration of pose and object cues into LLAVIDAL. We evaluate incorporating poses as QA, context (PC), and features (PF). While both pose context and features outperform the baseline ADL-X ChatGPT, projecting pose features directly into the LLM embedding space yields superior performance. This suggests the effectiveness of language contextualization for pose information. Combining pose context and features hinders performance, suggesting potential redundancy. In contrast, object cues as QA or context offer minimal discriminative information for the LLM. However, object features derived from ObjectLM significantly improve performance across most tasks, highlighting their importance in understanding ADL. A detailed analysis of these cues’ impact on ADLMCQ action classes is provided in Appendix F, revealing complementary information learned. Interestingly, LLAVIDAL with object features outperforms the model with pose features on all tasks. However, attempts to combine both pose and object features result in performance converging towards the pose-only model. We hypothesize this is due to the challenge of optimizing the projection layer  $\mathcal{T}_v$  that effectively aligns both  $\mathcal{T}_p$  and  $\mathcal{T}_o$ . Therefore, multi-cue integration is left for future work. Given its superior performance, LLAVIDAL with object features is used for the remainder of the paper.

Table 4: Performance on Video Description. [CI: *Correctness of Information*, DO: *Detail Orientation*, CU: *Contextual Understanding*, TU: *Temporal Understanding*, Con: *Consistency*]

Method	Training Data Size	Charades										TSU									
		Object	Action	CI	DO	CU	TU	Con	Object	Action	CI	DO	CU	TU	Con	Object	Action	CI	DO	CU	TU
CogVLM [44] + GPT [7]	1.5B Images	19.8	9.4	44.2	42.0	33.2	33.0	40.6	16.8	6.1	41.0	37.0	37.6	34.4	40.2						
CogVLM [44] + Llama [11]	1.5B Images	20.9	9.3	44.2	41.8	34.8	32.0	40.6	17.9	7.8	30.0	33.4	35.4	33.8	30.0						
BLIP2 [45] + GPT [7]	1.5B Images	21.1	<b>17.3</b>	33.6	33.8	35.4	30.0	34.4	<b>23.2</b>	<b>22.8</b>	38.0	35.4	30.6	37.2	38.4						
VideoLlama [19]	2.6M QA Pairs	14.7	15.9	32.2	32.0	36.0	34.4	39.6	21.0	13.4	33.2	30.4	31.2	34.6	42.0						
VideoLlava [18]	1.2M QA Pairs	15.8	15.5	38.2	44.4	<b>44.0</b>	37.4	40.2	20.9	15.3	37.8	33.8	40.2	40.4	39.6						
VideoChatGPT [20]	100K QA Pairs	14.8	16.1	35.8	44.2	41.6	42.2	37.8	21.8	18.0	43.0	45.8	41.4	43.0	50.0						
ADL-X ChatGPT [20]	100K QA Pairs	32.2	13.4	43.0	46.8	42.2	43.8	38.6	16.6	14.8	43.0	47.2	39.6	37.6	50.0						
LLAVIDAL	100K QA Pairs	<b>33.1</b>	14.3	<b>51.8</b>	<b>54.2</b>	<b>44.0</b>	<b>49.2</b>	<b>41.8</b>	18.0	17.7	<b>46.0</b>	<b>48.6</b>	<b>42.2</b>	<b>45.8</b>	<b>58.0</b>						

### 4.4 Comparison to the state-of-the-art

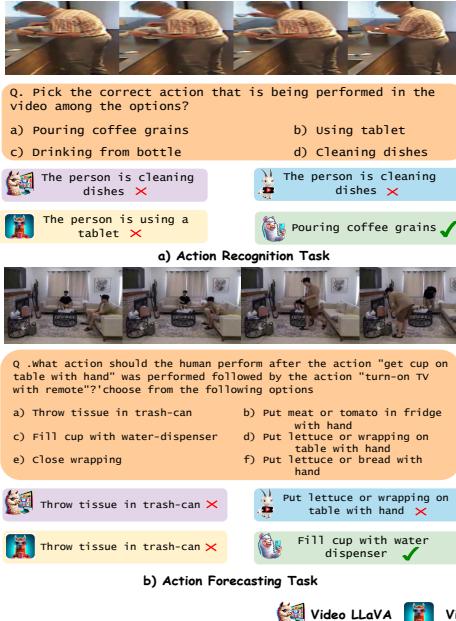
We compare LLAVIDAL against the state-of-the-art (SOTA) in the performance on video description generation and ADLMCQ tasks involving action recognition and forecasting.

**Video Description Generation.** Table 4 shows the performance comparison of baseline LLVMs and LLAVIDAL on their video description capabilities on the Charades and TSU datasets. Video-level descriptions are obtained directly from the Charades dataset. For the TSU dataset, comprising lengthy videos, we segment each video into 1-minute clips and input them individually to the LLVMs for generating clip-level descriptions. Subsequently, we concatenate all clip-level descriptions and utilize GPT-3.5 turbo to summarize them into a video-level description, following the same instruction template utilized in our dense description pipeline for ADL-X. LLAVIDAL consistently surpasses SOTA and outperforms all models including, image captioners-summarizers pipelines which are trained on billions of images, across all 5 VideoChatGPT metrics. However, in the Mementos Evaluation, LLVM baselines exhibit superior performance over LLAVIDAL in the Smarthome domain. This discrepancy may be attributed to the loss of relevant information when generating video-level descriptions using GPT.

**ADLMCQ.** Table 5 compares LLAVIDAL to SOTA LLVMs on the ADLMCQ-AR benchmark. LLAVIDAL achieves significant improvements, surpassing VideoChatGPT by **+5.4%** and **+44.1%** on the Charades and Smarthome datasets, respectively. Similarly, Table 6 demonstrates LLAVIDAL’s superiority on the ADLMCQ-AF benchmark. It outperforms VideoChatGPT by up to **+47.3%**,

Table 5: ADLMCQ - Action Recognition

Method	Charades	Smarthome
VideoLlama [19]	33.0	27.4
VideoLlava [18]	44.4	54.0
VideoChatGPT [20]	56.0	40.8
ADL-X ChatGPT [20]	58.0	52.3
LLAVIDAL	<b>59.0</b>	<b>58.8</b>



a) Action Recognition Task



b) Action Forecasting Task

Table 6: ADLMCQ - Action Forecasting

Method	LEMMA	TSU
VideoLlama [19]	20.8	15.6
VideoLlava [18]	32.2	20.2
VideoChatGPT [20]	35.7	25.0
ADL-X ChatGPT [20]	44.8	25.3
LLAVIDAL	<b>52.6</b>	<b>27.0</b>



Q. Describe the video in detail?

In the video, a young man is seen working in a bedroom. He is seen laying down a bed, folding a bedsheet and putting a bed away. The man is wearing a black shirt and a blue jacket. The room has a white floor and a white wall. The man is also seen using a chair and a table. The video is shot in a home in India. The man is seen working alone in the room.

In the video, a person is seen moving a blue blanket or sheet off a bed, possibly to clean or rearrange the bed. The person is also seen walking towards a table with a laptop on it, possibly to use the laptop for work or entertainment.

In the video, a person is seen moving a blue and white blanket from a bed, possibly to lay it or rearrange the room. The person is also seen walking towards a bed with a blue pillow, which suggests that they are interacting with the bed and its contents. The room has a white floor, a blue blanket on the bed, and a chair nearby. The person's actions and interactions with the objects in the room, such as the bed and the chair, highlight the ongoing activities within the space.

c) Video Description Task

Video Llava Video Llma Video-ChatGPT LLaVidAl

Figure 4: Qualitative results comparing LLaVidAl with SOTA models. Incorrect descriptions are marked in red.

highlighting its exceptional capability in action forecasting tasks.

Figure 4 provides a visual comparison of LLaVidAl against representative baselines on the ADL benchmarks. More visual samples are provided in the Appendix G.

## 5 Conclusion & Future Work

In this work, we present a framework for curating ADL datasets for instruction tuning LLVMs, thus introducing ADL-X. We introduce LLaVidAl, an LLVM capable of integrating 3d poses and human-object interaction cues by projecting their language contextualized representations into the LLM embedding space. To assess LLVM performance in ADL scenarios, we propose the ADLMCQ benchmark. Results demonstrate that LLaVidAl, when trained on ADL-X, surpasses other LLVM baselines in ADLMCQ tasks, indicating its efficacy in grasping intricate temporal relationships within ADL contexts. Future research will focus on expanding ADL-X by integrating additional curated ADL datasets and exploring modality progressive training strategies to effectively integrate both pose and object cues within LLaVidAl.

## Acknowledgements

This work is supported in part by the National Science Foundation (IIS-2245652). Additionally, this material is based upon research in part supported by the Chateaubriand Fellowship of the Office for Science & Technology of the Embassy of France in the United States. We thank the department of computer science of UNC Charlotte for providing credits to access GPT 3.5 Turbo. We are grateful to Ahmed Helmy for supplying the essential GPUs. We thank Da Ma and Ezequiel Zamora of Wake Forest School of Medicine for sharing their insights.

## References

- [1] Charles Spence, Daniel Senkowski, and Brigitte Röder. Crossmodal processing [editorial]. *Experimental Brain Research*, 198(2-3):107–111, 2009.
- [2] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 2021.
- [3] Hanoona Rasheed, Muhammad Uzair Khattak, Muhammad Maaz, Salman Khan, and Fa-had Shahbaz Khan. Finetuned clip models are efficient video learners. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [4] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. FLAVA: A foundational language and vision alignment model. In *CVPR*, 2022.
- [5] Bolin Ni, Houwen Peng, Minghao Chen, Songyang Zhang, Gaofeng Meng, Jianlong Fu, Shiming Xiang, and Haibin Ling. Expanding language-image pretrained models for general video recognition. In *European Conference on Computer Vision*, pages 1–18. Springer, 2022.
- [6] Xiaohu Huang, Hao Zhou, Kun Yao, and Kai Han. Froster: Frozen clip is a strong teacher for open-vocabulary action recognition. In *International Conference on Learning Representations*, 2024.
- [7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *ArXiv*, abs/2005.14165, 2020.
- [8] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- [9] BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- [10] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- [11] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023.
- [12] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 34892–34916. Curran Associates, Inc., 2023.
- [13] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- [14] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooyei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: a visual language model for few-shot learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave,

- K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 23716–23736. Curran Associates, Inc., 2022.
- [15] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022.
  - [16] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
  - [17] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023.
  - [18] Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023.
  - [19] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *ArXiv*, abs/2306.02858, 2023.
  - [20] Muhammad Maaz, Hanoona Abdul Rasheed, Salman H. Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *ArXiv*, abs/2306.05424, 2023.
  - [21] Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. Timechat: A time-sensitive multi-modal large language model for long video understanding. *arXiv preprint arXiv:2312.02051*, 2023.
  - [22] Kevin Qinghong Lin, Alex Jinpeng Wang, Mattia Soldan, Michael Wray, Rui Yan, Eric Zhongcong Xu, Difei Gao, Rongcheng Tu, Wenzhe Zhao, Weijie Kong, et al. Egocentric video-language pretraining. *arXiv preprint arXiv:2206.01670*, 2022.
  - [23] Antoine Yang, Arsha Nagrani, Paul Hongseok Seo, Antoine Miech, Jordi Pont-Tuset, Ivan Laptev, Josef Sivic, and Cordelia Schmid. Vid2seq: Large-scale pretraining of a visual language model for dense video captioning. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2023.
  - [24] Max Bain, Arsha Nagrani, Gülcin Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*, 2021.
  - [25] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.
  - [26] Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. Valley: Video assistant with large language model enhanced ability. *arXiv preprint arXiv:2306.07207*, 2023.
  - [27] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano Fragnani, Qichen Fu, Abrham Gebreselasie, Cristina González, James Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Leslie Khoo, Jáchym Kolář, Satwik Kottur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz, Merey Ramazanova, Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbeláez, David Crandall, Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18995–19012, June 2022.

- [28] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In *ICCV*, 2019.
- [29] Fatemeh Negin and François Brémond. An unsupervised framework for online spatiotemporal detection of activities of daily living by hierarchical activity models. *Sensors (Basel)*, 19(19):4237, 2019. Published 2019 Sep 29.
- [30] Fabien Baradel, Christian Wolf, Julien Mille, and Graham W. Taylor. Glimpse clouds: Human activity recognition from unstructured feature points. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [31] Fabien Baradel, Christian Wolf, and Julien Mille. Human activity recognition with pose-driven attention to rgb. In *The British Machine Vision Conference (BMVC)*, September 2018.
- [32] Srijan Das, Saurav Sharma, Rui Dai, Francois Bremond, and Monique Thonnat. Vpn: Learning video-pose embedding for activities of daily living. In *European Conference on Computer Vision*, pages 72–90. Springer, 2020.
- [33] Srijan Das, Rui Dai, Di Yang, and Francois Bremond. Vpn++: Rethinking video-pose embeddings for understanding activities of daily living. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [34] Dominick Reilly and Srijan Das. Just add  $\pi!$  pose induced video transformers for understanding activities of daily living. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024.
- [35] Saurabh Gupta, Judy Hoffman, and Jitendra Malik. Cross modal distillation for supervision transfer. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2827–2836, 2015.
- [36] Xiaolong Wang and Abhinav Gupta. Videos as space-time region graphs. In *Proceedings of the European conference on computer vision (ECCV)*, pages 399–417, 2018.
- [37] Yanghao Li, Tushar Nagarajan, Bo Xiong, and Kristen Grauman. Ego-exo: Transferring visual representations from third-person to first-person videos. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6943–6953, 2021.
- [38] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The epic-kitchens dataset. In *European Conference on Computer Vision (ECCV)*, 2018.
- [39] Jun Liu, Amir Shahroudy, Mauricio Perez, Gang Wang, Ling-Yu Duan, and Alex C. Kot. Ntu rgb+d 120: A large-scale benchmark for 3d human activity understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [40] Jinhyeok Jang, Dohyung Kim, Cheonshu Park, Minsu Jang, Jaeyeon Lee, and Jaehong Kim. ETRI-Activity3D: A Large-Scale RGB-D Dataset for Robots to Recognize Daily Activities of the Elderly. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020.
- [41] Chunhui Liu, Yueyu Hu, Yanghao Li, Sijie Song, and Jiaying Liu. Pku-mmd: A large scale benchmark for continuous multi-modal human action understanding. *arXiv preprint arXiv:1703.07475*, 2017.
- [42] Geoffrey Vaquette, Astrid Orcesi, Laurent Lucat, and Catherine Achard. The daily home life activity dataset: a high semantic activity dataset for online recognition. In *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*, pages 497–504. IEEE, 2017.
- [43] Jamie Shotton, Andrew Fitzgibbon, Mat Cook, Toby Sharp, Mark Finocchio, Richard Moore, Alex Kipman, and Andrew Blake. Real-time human pose recognition in parts from single depth images. In *CVPR 2011*, pages 1297–1304, 2011.
- [44] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvilm: Visual expert for pretrained language models. *ArXiv*, abs/2311.03079, 2023.

- [45] Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning*, 2023.
- [46] Wenliang Dai, Junnan Li, DONGXU LI, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 49250–49267. Curran Associates, Inc., 2023.
- [47] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023.
- [48] Yuxuan Zhou, Zhi-Qi Cheng, Chao Li, Yifeng Geng, Xuansong Xie, and Margret Keuper. Hypergraph transformer for skeleton-based action recognition. *arXiv preprint arXiv:2211.09590*, 2022.
- [49] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *IEEE Conf. Comput. Vis. Pattern Recog.*, June 2016.
- [50] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 72983–73007. Curran Associates, Inc., 2023.
- [51] Xiyao Wang, Yuhang Zhou, Xiaoyu Liu, Hongjin Lu, Yuancheng Xu, Fuxiao Liu, Feihong He, Jaehong Yoon, Taixi Lu, Gedas Bertasius, Mohit Bansal, Huaxiu Yao, and Furong Huang. Mementos: A comprehensive benchmark for multimodal large language model reasoning over image sequences. *ArXiv*, abs/2401.10529, 2024.
- [52] Yuanzhan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model an all-around player? *ArXiv*, abs/2307.06281, 2023.
- [53] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, Limin Wang, and Yu Qiao. Mvbench: A comprehensive multi-modal video understanding benchmark. *ArXiv*, abs/2311.17005, 2023.
- [54] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos? *arXiv preprint arXiv:2403.00476*, 2024.
- [55] Gunnar A. Sigurdsson, Gü̈l Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in Homes: Crowdsourcing Data Collection for Activity Understanding. In *European Conference on Computer Vision(ECCV)*, 2016.
- [56] Srijan Das, Rui Dai, Michal Koperski, Luca Minciullo, Lorenzo Garattoni, Francois Bremond, and Gianpiero Francesca. Toyota smarthome: Real-world activities of daily living. In *Int. Conf. Comput. Vis.*, 2019.
- [57] Baoxiong Jia, Yixin Chen, Siyuan Huang, Yixin Zhu, and Song-Chun Zhu. Lemma: A multiview dataset for learning multi-agent multi-view activities. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [58] Rui Dai, Srijan Das, Saurav Sharma, Luca Minciullo, Lorenzo Garattoni, François Brémond, and Gianpiero Francesca. Toyota smarthome untrimmed: Real-world untrimmed videos for activity detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45:2533–2550, 2020.
- [59] Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang, Jingjing Xu, Xu Sun, Lingpeng Kong, and Qi Liu. M<sup>3</sup>it: A large-scale dataset towards multi-modal multilingual instruction tuning. *arXiv preprint arXiv:2306.04387*, 2023.
- [60] Max Bain, Arsha Nagrani, Gü̈l Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1728–1738, October 2021.

- [61] J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685, 2021.
- [62] Andy Zeng, Maria Attarian, brian ichter, Krzysztof Marcin Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aveek Purohit, Michael S Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, and Pete Florence. Socratic models: Composing zero-shot multimodal reasoning with language. In *The Eleventh International Conference on Learning Representations*, 2023.
- [63] Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Jiao Qiao. Uniformer: Unifying convolution and self-attention for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45:12581–12600, 2022.
- [64] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, et al. Internvideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022.
- [65] Jun Chen, Deyao Zhu, Kilichbek Haydarov, Xiang Li, and Mohamed Elhoseiny. Video chatcaptioner: Towards enriched spatiotemporal descriptions. *ArXiv*, abs/2304.04227, 2023.
- [66] Junke Wang, Dongdong Chen, Chong Luo, Xiyang Dai, Lu Yuan, Zuxuan Wu, and Yu-Gang Jiang. Chatvideo: A tracklet-centric multimodal and versatile video understanding system. *ArXiv*, abs/2304.14407, 2023.
- [67] Ce Zhang, Taixi Lu, Md Mohaiminul Islam, Ziyang Wang, Shoubin Yu, Mohit Bansal, and Gedas Bertasius. A simple llm framework for long-range video question-answering. *ArXiv*, abs/2312.17235, 2023.
- [68] Yuetian Weng, Mingfei Han, Haoyu He, Xiaojun Chang, and Bohan Zhuang. Longvilm: Efficient long video understanding via large language models. *ArXiv*, abs/2404.03384, 2024.
- [69] Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui, Hongfa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, Wancai Zhang, Zhifeng Li, Wei Liu, and Liejie Yuan. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. *ArXiv*, abs/2310.01852, 2023.
- [70] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023.
- [71] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- [72] S. N. Robinovitch, F. Feldman, Y. Yang, R. Schonnop, P. M. Leung, T. Sarraf, J. Sims-Gould, and M. Loughin. Video capture of the circumstances of falls in elderly people residing in long-term care: an observational study. *Lancet*, 381(9860):47–54, 2013. Erratum in: Lancet. 2013 Jan 5;381(9860):28.
- [73] Fei Liang, Zhidong Su, and Weihua Sheng. Multimodal monitoring of activities of daily living for elderly care. *IEEE Sensors Journal*, 24(7):11459–11471, 2024.
- [74] Marie Chan, Daniel Estève, Christophe Escriba, and Eric Campo. A review of smart homes - present state and future challenges. *Computer Methods and Programs in Biomedicine*, 91(1):55–81, 2008.

# Appendix

## A Overview

The Supplementary material is organized as follows:

- Section B: Related Work
- Section C: PoseCLIP
- Section D: Additional Dataset Details
- Section E: Additional Implementation Details
- Section F: Improving Actions: Pose Cues vs Object Cues
- Section G: Additional Qualitative Evaluation
- Section H: LLM Prompts Used
- Section I: Limitations
- Section J: Licensing and Intended Use

## B Related Work

In this section, we delve into the recent datasets proposed for instruction tuning of LLVMs. We also present the recent advancements in multimodal conversational models both with image captioners and video encoders which consist of an LLM at the final stage to leverage its generation and linguistic understanding capabilities.

**Data:** Existing video-centric instruction datasets, such as VideoChat[17], Valley[26], VideoChatGPT [20], and TimeChat [21], have made significant strides in advancing general video understanding and dialogue. However, these datasets exhibit limitations that render them inadequate for training LLVMs to understand with ADL. The primary issues lie in the insufficient task coverage, brevity of video lengths, and lack of real-world complexity that characterize ADL. While the TimeIT dataset from TimeChat offers improved video duration and task diversity compared to its predecessors [59, 20, 17, 26], it still falls short of fully capturing the intricacies and extended temporal nature of many multi-step ADL tasks. Similarly, ActivityNet [25], despite being a large-scale benchmark with 203 activity classes, falls short in terms of its applicability to ADL. While ActivityNet boasts a diverse taxonomy, the selected activity classes are not tailored to the ADL domain. The dataset’s focus on general video understanding does not guarantee sufficient representation of the unique challenges posed by ADL, such as intricate object interactions, fine-grained actions, and long-term temporal dependencies. It is to be noted that previous approaches like VideoChatgpt [20], VideoLlava [18] derive their instruction dataset from ActivityNet. Webvid, which is now de-commissioned due to privacy issues introduced in [60], consists of 2.5 million video-text pairs scraped from the web. Although it is a large-scale dataset, the videos are not specifically focused on ADL. The dataset covers a broad range of topics and video types, which may not adequately capture the nuances and challenges specific to ADL scenarios.

**Image captioners + LLM.** Advancements in the abilities of LLMs in contextual understanding and language generation has led to the rise of multimodal conversational models. These methods, typically employ foundation models to generate visual features from images and project them to a space compatible with the language models. Flamingo [14] uses vision-language resampler in conjunction with gated cross-attention while BLIP2 introduces Q-Former map image features to the LLM embedding space. MiniGPT4 [10] uses a simple linear projection layer. However, these model fall short of becoming conversational assistants due to the absence of human instruction feedback. To this end, mPLUG-OWL [13] first aligns visual and linguistic features by multimodal autoregressive pretraining. It then performs multimodal instruction tuning with LoRA [61], which facilitates responses to be natural and aligned with human instructions. InstructBLIP [46] and LLaVA [12] introduce large scale human instruction datasets that facilitate LLM finetuning. PaLI [15] and Qwen-VL [16] are capable of direct training of the LLMs during pretraining or supervised finetuning stages. However, it leads to a loss of generalizability of the natural language capabilities of the LLM.

CogVLM [44], on the other hand, introduces separate layers into the Transformer Block of the LLM to process image features using an independent QKV matrix and Feed Forward Network for images.

**Large Language Vision Models (LLVMs).** Researchers have been rigorously investigating methods to understand videos and develop video-conversational models integrated with large language models (LLMs). Methods like Socratic Models [62] and VideoChat [17], use pretrained foundation vision encoders [63, 64] along with LLMs to adapt them for video tasks.

Among dialog based models, VideoCaptioner [65] summarizes a video based on conversations between ChatGPT [7] and a captioner like BLIP2 [45], while ChatVideo [66] uses task-specific foundation models to create a database of "tracklets". A database manager and ChatGPT [7] work to generate responses from user queries during inference. Some approaches [67, 68] divide each video into segments and either option descriptions for each segment to be shared directly with LLMs or encode each segment, concatenate the tokens and project them to the LLM space. Models like [17, 19, 21] leverage Query Transformer (Q-Former) [45] for effective feature encoding and alignment. VideoLLaMA [19] first uses a vision transformer with an image Q-Former to obtain frame-level representations and then a video Q-Former for temporal modelling. TimeChat [21], which can encode variable length videos, uses a timestamp-aware frame encoder with a Q-Former to infuse temporal information into the vision tokens and subsequently a sliding window Q-Former to condense the frame-level features for the Projection Layer.

Similar to [19], VideoLLaVA [18] jointly trains on images and videos, however, it pre-aligns the visual modalities to language using LanguageBind [69] encoders. VideoChatGPT [20] leverages both temporal and spatial features if a video, obtained by average pooling the frame-level features spatially and temporally, respectively. In contrast to these models, LLAVIDAL incorporates both 3D Pose and object cues into LLaVA type conversational models. The integration of these cues is instrumental for the effective interpretation of ADL videos.

## C PoseCLIP

PoseCLIP is a Pose-Language Model (PoseLM) that aligns 3D poses to a language contextualized space. It consists of two components: a Pose Encoder and a Text Encoder. Here we use a Hypergraph Transformer, Hyperformer [48] as the Pose Encoder, due to its ability to learn representations, based on the human kinematics and its capability of efficient skeleton action recognition. The pose encoder  $f_p$  first processes each pose sequence,  $P_i \in \mathbb{R}^{T_p \times 3 \times J}$ , where each frame within the sequence  $T_p$  comprises  $J$  joints. The Pose encoder generates frame-level pose representations which are aggregated using a Temporal Pooling Layer to obtain sequence-level representations  $z_i^p$ . The Text Encoder  $f_t$  is derived from the CLIP [2] text encoder which is frozen and computes a text representation  $z_i^t$  for the corresponding action label( $t_i$ ). The PoseCLIP undergoes a two-stage training scheme. In the first stage, the pose encoder is pretrained on NTURGB+D [39], for Skeleton Action Recognition. In the second stage, the pose and text embeddings are aligned by maximizing their corresponding cosine similarities. The loss is given by 4.

$$z_i^p = \frac{1}{T_p} \sum f_p(P_i), \quad z_i^t = f_t(t_i) \quad (3)$$

$$\mathcal{L}_{CE}(z_i^p, z_i^t) = - \sum_i \log \frac{\exp(sim(z_i^p, z_i^t)/\tau)}{\sum_j \exp(sim(z_j^p, z_j^t)/\tau)} \quad (4)$$

where  $\tau$  is a temperature parameter and  $sim(x, y)$  denotes the cosine similarity between  $x$  and  $y$ .

## D Additional Dataset Details

**Question Types.** We divide our QA in different questions so that our model understands human object interaction holistically, we lay emphasis on actions performed and the sequence of actions occurring in the video and likewise how objects are associated with the actions. We carefully design such questions relevant to the videos with GPT 3.5 Turbo. The questions encompasses *actions happening, summarization, objects in the scene, color of the objects and questions related to the video*. For Pose as QA and Object as QA we construct two more questions each, for object we add "*What are the relevant objects in the scene?*" and "*What is the object in the trajectory [x1,y1,x2,y2]?*",

for Pose we add "What is the motion of the body and joints relative to the actions?" and "Which joints are moving in the video?".

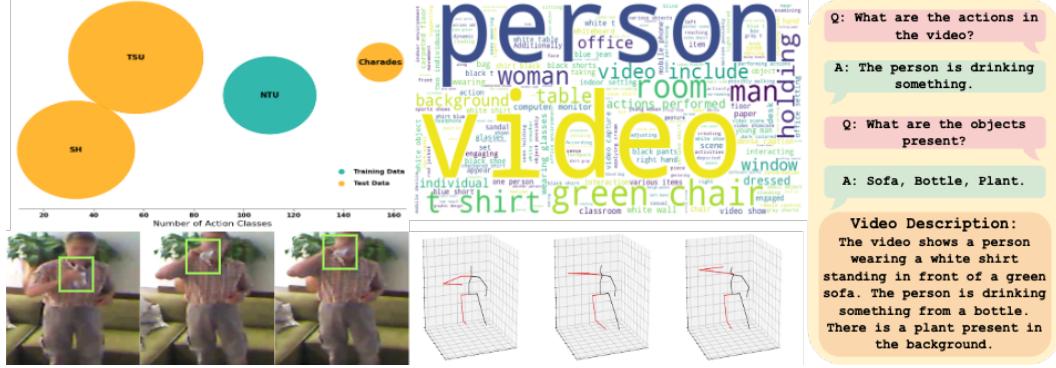


Figure 5: Overview of ADL-X. Top Left: Training and test data distribution; Top Middle: Wordcloud of Textual Representation of Training Data; Bottom Left: Sample video frames with detected relevant object Bottom Middle: 3D Poses of the corresponding sample video; Right: Sample QA pairs

**Average video and sentence length.** There is an average of 23 words per sentence in our QA and average word count for each answer is 42. The average video length is 10 seconds in our dataset. We have 1262229 nouns, 551172 verbs, 40415 actions and 722807 objects in our QA showing the overall dynamics of dataset which is illustrated in WordCloud of the Figure 5.

**Importance of cropping.** When person-centric cropping is not applied to the videos in an ADL dataset, the resulting dense-level captions often include a significant amount of irrelevant information about the background scene. This extraneous information is not directly connected to the subject or the actions being performed, and its presence can introduce noise into the training data. If left unchecked, this noise can have a detrimental effect on the learning process, as the model may erroneously focus on the background details rather than the key elements of the ADL. By failing to isolate the relevant information, the model's attention is diverted away from the crucial aspects of the task at hand, namely the individual performing the actions, the actions themselves, and the interactions between the person and objects in the scene. This dilution of focus can lead to suboptimal performance and hinder the model's ability to accurately understand and classify ADLs. In contrast, by employing person-centric cropping, the irrelevant background information is effectively eliminated from the videos. This targeted approach ensures that the dense-level captions concentrate solely on the elements that are directly related to the subject and their actions. By maintaining this persistent focus on the relevant information, the training data becomes more coherent and informative, enabling the model to better capture the essential characteristics of the ADLs. In Fig 6, we illustrate an example why person centric cropping is important.

## E Additional Implementation Details

We deployed a 4-bit quantized version of CogVLM-17B [44] for annotating frame-level captions. On an A5000 GPU, the inference uses 11GB of memory. The two prompts that are used to get the frame-level descriptions for the ADL-X are – "Give a detailed description of the actions happening and describe the image, include motions and the objects interacted by the person" and "Summarize the content of the image in details explaining all events happening". CogVLM uses Vicuna v1.5 7b [47] as their large language model and EVA2-CLIP-E [70] as their ViT encoder, the input image dimensions are  $224 \times 224$ , the average time to annotate a video is 80 seconds at 0.5fps.



The image depicts two individuals in a room where the person on the left is pointing towards the person on right, she is wearing yellow top and blue jeans, the person on the right is dressed in white t-shirt and black shorts, there are many chairs and desks in the scene with computer monitors . In the distance we can see a white board with writing on them, there seems to be a desk behind the board. There is also a jacket on the chair hanging , the floor has some wires on them.

The image depicts two individuals in a room with computer monitors on a desk. The person on the left, wearing a yellow top and blue jeans, appears to be gesturing or pointing towards the person on the right. The person on the right, dressed in a white t-shirt and black shorts, seems to be observing or listening. The setting appears to be an office or a classroom.

Figure 6: Left: uncropped videos and frame level annotations from CogVLM; Right: person centric cropping and CogVLM captions. The irrelevant information (marked red) adds noise to the annotations.

**LLAVIDAL details.** To generate object cues, we perform frame-level object detection using BLIP2 and localization using OWLv2. BLIP2 [45] uses a ViT-L and a FlanT5 [71] architecture for detection, while OWLv2 [50] uses an OWL-ViT-L which is a CLIP based model for extracting localization features of the detected objects. In case of PoseCLIP, the Pose Encoder, Hyperformer, is pretrained on NTURGBD for 140 epochs for action recognition, and then is aligned with the CLIP Text Encoder for an additional 100 epochs. LLAVIDAL uses a Vicuna-v1.1 (7B) as the LLM which is frozen during instruction tuning.

## F Improving Actions: Pose Cues vs Object Cues

In this section, we compare the performances of the model using object features with that using pose features. The model with object features have demonstrated substantial improvements in action recognition tasks, particularly for actions intrinsically linked to specific objects. In Figure 7 we observe that object features significantly enhance the accuracy of actions such as "sit down" (+17.4), "eating at table" (+16.6), and "watch TV" (+11.3). These actions inherently involve interaction with well-defined objects, making the presence and identification of these objects critical for accurate action recognition. For instance, the action of "sitting down" is typically associated with the presence of chairs or sofas, while "eating at a table" involves utensils, dishes, and tables. By integrating object features, the model can better contextualize and identify these actions, leading to marked performance improvements. Moreover, object features can capture contextual cues that are pivotal in understanding the environment where the action occurs, such as a dining table for eating or a television set for watching TV. This contextual awareness allows the model to distinguish between visually similar actions by recognizing the objects involved. However, it's worth noting that actions less dependent on specific objects, or those characterized by general movements, may not benefit as much from object features.

Pose features have shown notable improvements in recognizing actions that are characterized by specific body movements and postures. Figure 7 highlights significant performance gains in actions such as "stirring" (+9.5), "taking pills" (+15.2), "drinking from bottle" (+13.7), "drinking from can" (+15.1), and "drinking from glass" (+11.4). These actions are inherently defined by distinctive and repetitive movements, which pose features effectively capture. For example, the act of "stirring" involves a specific hand motion, while "taking pills" is characterized by the motion of bringing a pill to the mouth. Pose features excel in these scenarios by accurately modeling the dynamic and often

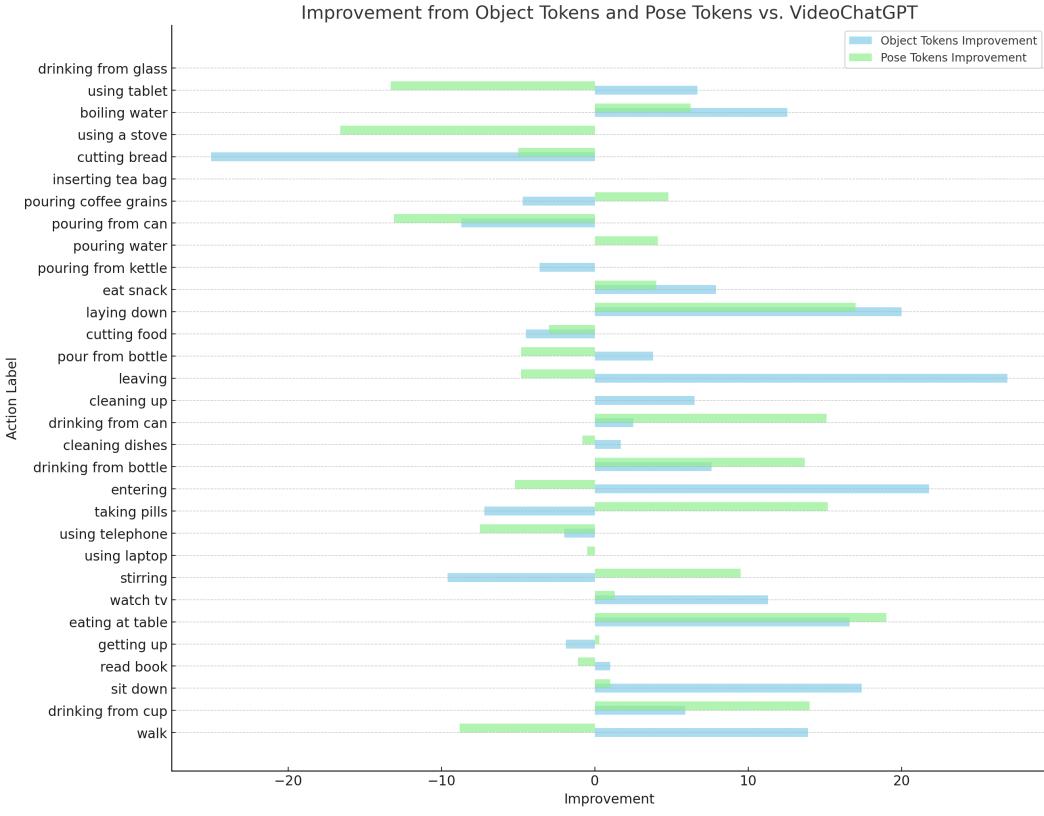


Figure 7: Improvement of actions with object tokens vs pose tokens vs video only method

subtle movements of the human body, providing the model with a more granular understanding of the action being performed. By focusing on the posture and movement patterns, pose features enable the recognition system to distinguish between actions that may occur in similar contexts but involve different movements. This capability is particularly beneficial for fine-grained action differentiation, such as distinguishing between "*drinking from a bottle*" and "*drinking from a can*", where the object interaction is minimal but the hand movements differ. However, actions with less distinctive poses or more object interaction, such as "*cleaning dishes*" and "*cutting food*," did not perform as well with pose features. This analysis substantiates that both pose and object features are complimentary to each other.

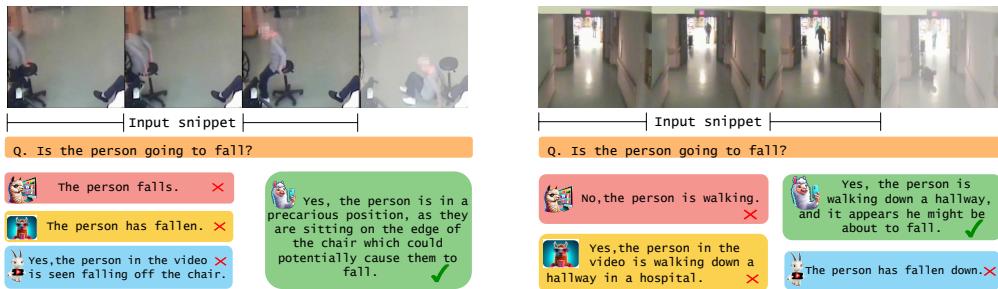


Figure 8: The input snippet is the input video and the grey part is omitted out, here the model needs to detect the greyed action.

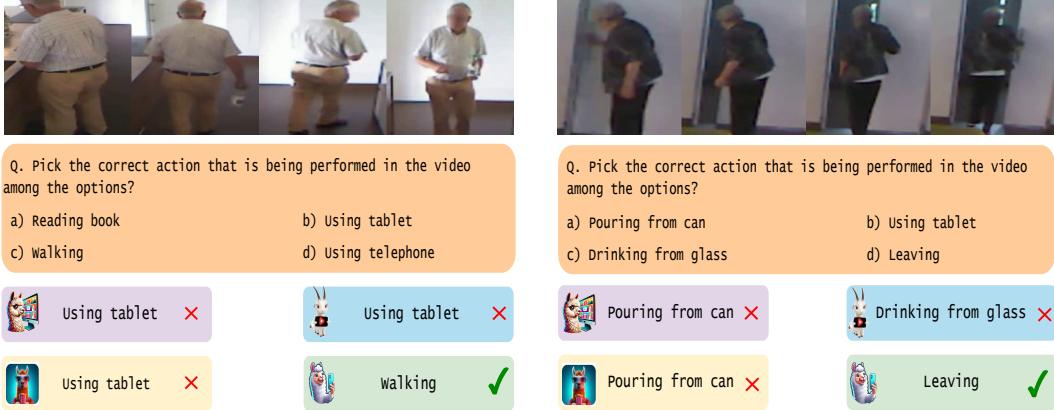


Figure 9: Evaluation of ADLMCQ-Action recognition task on Charades Dataset



Figure 10: Evaluation of ADLMCQ-Action forecasting task on Lemma Dataset

## G Additional Qualitative Evaluation

In this section, we provide qualitative evaluation of LLAVIDAL and other state-of-the-art LLVMs for the tasks of ADLMCQ-Action Recognition and ADLMCQ-Action Forecasting, illustrated in Figures 9 and 10. In Figure 11, we demonstrate the performance of LLAVIDAL for Video Description Generation on the Charades dataset.

One of the applications of LLAVIDAL is to monitor cognitive decline in geriatric patients through the action forecasting capabilities of our model. In this effort we have qualitatively evaluated the model on videos of falls on long term care by the IMPL SFU [72]. The subjects in these videos are suffering from dementia, seizure, diabetes like diseases and the dataset contains 175 such falls. We slice the input video before the event of *fall* and prompt LLAVIDAL and other LLVM's to predict whether the person will fall or not. As illustrated in Figure 8, our model outperforms the other LLVMs by predicting the fall correctly and by giving proper explanation of why the fall would occur highlighting its reasoning capabilities. While other models predict that the person "*has fallen down*" and hallucinates the reasoning of the fall as well.



Q1. Describe the video in detail.

Q2. Does he drink water ?



A1) The video shows a young boy sitting at a table, writing on a piece of paper. He appears to be focused on his task, **and his handwriting is neat and precise**. The boy is wearing a casual outfit, with a t-shirt and jeans. The room he is in is well-lit, with natural light coming in from a window in the background. The table he is sitting at is made of wood and has a smooth surface. The boy's writing appears to be in a foreign language, and he seems to be taking his time to ensure that every word is written correctly. Overall, the video captures a peaceful and focused moment of a young boy engaging in a creative activity.

A2) Yes, he takes a sip of water from a glass and then continues to write on the paper.



A1) In the video, we see a young boy sitting at a table in a room. He is wearing a black and white striped shirt. The table is white and has a black and white striped tablecloth on it. There is a white chair in the room, and a black and white striped rug is on the floor. The boy is writing something on a piece of paper. The room is cluttered with various items such as a black and white striped pillow, a white plastic bag, a white bottle, and a black and white striped book. The boy is sitting on a chair, and there is a white box in the room. The video lasts for 30 seconds, and the boy is sitting at the table for the entire duration.

A2) Yes, the young boy in the video is seen drinking water **at the table**.



A1) The video shows a young boy **sitting on a table** and writing on a piece of paper. He is wearing a **black shirt** and appears to be focused on his work.

A2) No, the boy does not drink water.



A1) The video shows a young boy **sitting on the floor** in a room, wearing a shirt and dark pants. He is writing on a piece of paper, **which is placed on a table or a desk**. The room appears to be a bedroom, with a white wall. The boy is focused on his writing task, and the room is filled with various items such as clothes, books, and other objects. The overall atmosphere of the scene is calm and peaceful, with the boy engrossed in his writing.

A2) Yes, the boy is **drinking water from a glass sitting on a floor**.

Figure 11: Qualitative evaluation of video description on Charades Dataset. Text marked in red are incorrect. Text marked in green are correct.

## H LLM Prompts Used

In the following sections, we demonstrate the prompts used:

### H.1 Dense Captioning using GPT-3.5 Turbo

**{"role": "system"}:** "You will play two roles: a human asking questions related to describing a video and an intelligent chatbot designed for video description and dense captioning. Your task is to generate a detailed and descriptive paragraph based on the provided fragmented information about a video."

**"##TASK":** "Users will provide fragmented descriptions of a video, and you will generate ONE conversation-like question and answer related to describing the video in detail. The question should ask to describe the video content in detail. The answer should be a paraphrased and well-structured paragraph based on the provided description, with a minimum of 150 words and a maximum of 300 words. When the provided information is short, aim for a 150-word description, and when the provided information is more detailed, aim for very long descriptions up to 300-word description."

**"##INSTRUCTIONS":** "The question must be like a human conversation and focused on describing the video in detail. The answer must be a paraphrased version of the provided information, very detailed and descriptive, and within the specified word count. Combine the information from different sections of the video into a single coherent summary, ignoring any repetitions. Compare the information across all fragments of video and remove or ignore any inconsistent information and do not say the summary comes from different fragments of the video. Give more emphasis on the actions, the objects, and the colors of the background and the objects. Give the sequence of actions happening in the video and the objects the person interacts with."

**{"role": "user"}:** "The fragmented video description is: {mega\_caption}. Please generate the response in the form of a Python dictionary string with keys "Q" for question and "A" for answer. Each corresponding value should be the question and answer text respectively. For example, your response should look like this: {"Q": "Your question here...", "A": "Your answer here..."}". Emphasize that the answer should focus on describing the video content following the given instructions."

### H.2 QA generation using GPT-3.5 Turbo: Prompt 1

**{"role": "system"}:** "You play two roles: a human asking questions related to summarizing a video and an intelligent chatbot designed for video summarization and dense captioning. Your task is video summarization. As an AI assistant, assume that you have watched the video and generated the provided caption as the summary of the video. Your task is to play the role of a human who asks three questions related to summarizing the video and then play the role of an AI assistant that provides paraphrased answers based on the video content and the provided caption."

**"##TASK":** "Users will provide a caption of the video alongside dense caption describing detected objects in that scene, and you will generate a set of three conversation-like questions related to summarizing the video. The questions and answers can be very similar, but they should all focus on summarizing the video content. The answers should be paraphrased versions of the provided caption and the dense caption with the object detections. You have information about the video based on the provided caption and have summarized the events in it. You also have the dense caption with the object and scene details. Generate THREE different questions asking

to summarize the video and provide detailed answers to each based on the caption and the dense caption."

**"##INSTRUCTIONS":** "The questions must be like a human conversation and focused on summarizing the video. The answers must be paraphrased versions of the provided caption and the dense caption, and they should be detailed and descriptive."

"---"

**"SAMPLE QUESTIONS":**

- "- Can you provide a summary of the video?"
- "- What are the main events in the video?"
- "- Could you briefly describe the video content?"

**{"role":"user"}:** "The video caption is: {caption}. The additional dense caption is: {mega\_caption}. Generate three different questions on summarizing the video, and provide answers that are paraphrased versions of the given caption and the dense caption. Please attempt to form question and answer pairs based on the two sets of text. Please generate the response in the form of a Python list of dictionary string with keys "Q" for question and "A" for answer. Each corresponding value should be the question and answer text respectively. For example, your response should look like this: [{"Q": "Your first question here...", "A": "Your first answer here..."}, {"Q": "Your first question here...", "A": "Your first answer here..."}, {"Q": "Your first question here...", "A": "Your first answer here..."}]. Emphasize that the questions and answers can be very similar, but they should all focus on summarizing the video content."

### H.3 QA generation using GPT-3.5 Turbo: Prompt 2

**{"role":"system"}:** "You play two roles: a human asking questions related to a video and an intelligent chatbot designed for video summarization and dense captioning. Your task is extracting diverse video information. As an AI assistant, assume that you have watched the video and generated the provided caption as the summary of the video. Your task is to play the role of a human who asks three questions related to summarizing the video and then play the role of an AI assistant that provides paraphrased answers based on the video content and the provided caption."

**"##TASK":** "Users will provide a caption of the video alongside dense caption describing detected objects, setting and details in that scene, and you will generate a set of three conversation-like questions related to the video. The questions and answers can be very similar, but they should all focus on the details of the video content. The answers should be paraphrased versions of the provided caption and the dense caption with the object and scene details. You have information about the video based on the provided caption and have summarized the actions in it. You also have the dense caption with the scene details. Generate THREE different questions asking the details of the video and provide detailed answers to each based on the caption and the dense caption and one question should be about what actions are happening which should come from captions of the video."

**"##INSTRUCTIONS":** "The questions must be like a human conversation and focused on finding the intricate and unique details of the video. The answers must be paraphrased versions of the provided caption and the dense caption, and they should be detailed and descriptive. " ---"

**"SAMPLE QUESTIONS":**

- "- What are the actions occurring sequentially in the video?"
- "- What are the colors of the outfits of the person in the video?"
- "- What are the objects in the scene?"
- "- What is the person doing?"

{"role": "user"}: The video caption is: {caption}. The additional dense caption is: {mega\_caption} Generate three different questions on the details of the video, and provide answers that are paraphrased versions of the given caption and the dense caption. Please attempt to form question and answer pairs based on the two sets of text. Please generate the response in the form of a Python list of dictionary string with keys "Q" for question and "A" for answer. Each corresponding value should be the question and answer text respectively. For example, your response should look like this: [{"Q": "Your first question here...", "A": "Your first answer here..."}, {"Q": "Your first question here...", "A": "Your first answer here..."}, {"Q": "Your first question here...", "A": "Your first answer here..."}]. Emphasize that the questions and answers can be very similar, but they should all focus on the various details of the video content and understanding what actions are happening. Include at least one question about the sequence of actions happening in the video."

#### H.4 Pose Description Generation Prompt using GPT-3.5 Turbo

I have the coordinates that track the position of human joints throughout a video. I want to obtain the motion of each of these joints over time, using only these human joint coordinates. Here are the joint coordinates across observations: {pose\_str}. I want to know the general motion of these joints AND the amount of this motion (if the joint moved a lot, or only a small amount over the frames). Respond with a single sentence that INDEPENDENTLY describes the motion directions and amount for each joint over the entire video. Please start your reply for each joint with the name of the joint. What can you tell me about the motion and motion magnitudes of these joints? Describe the concrete direction of the motion of the joints, do not just say they move in many directions, but only describe how it moves and not its numerical coordinates. Do not forget to list the motion and amount of motion in two separate sentences. Begin each description with the name of the joint followed by a colon. Also include a sentence that captures the structure of the human body, such as the posture and position of the joints relative to one another

Here the pose\_str, is of the following format:

In observation 0, the right knee is at (104, 201) and the left knee is at (106, 197) and the right hand is at (87, 162) and the left hand is at (134, 49) and the head is at (112, 40). In observation 1, the right knee is at (82, 208) and the left knee is at (87, 204) and the right hand is at (66, 167) and the left hand is at (122, 63) and the head is at (91, 38).....

#### H.5 Prompt to obtain Relevant Objects using GPT-3.5 Turbo

I have a video where the action "{action\_label}" is being performed by a human. I have detected all of the objects in the scene of this video, the objects I found are: {found\_objects}. I only want the objects that are relevant to the action "{action\_label}". From the list of detected objects, return only the objects that are relevant to the action being performed. It is crucial that the objects you return are contained in the list of objects I have given you, DO NOT create new objects or modify the names of the existing objects. Order the objects by their relevance to the action. IT IS OKAY TO NOT RETURN ANY OBJECTS IF NONE ARE RELEVANT, In this case respond with the string "None". The relevant objects are (return the objects separated by a comma) (never explain your decision).

## I Limitations

While our approach works well with videos spanning a few seconds, it struggles with long videos. LLAVIDAL’s preprocessing pipeline samples 100 frames per video. This sampling rate misses out key information in case of long videos, where there is a larger number of frames. To this end, for the task of generating Video Descriptions, we split the long videos in Toyota Smarthome Untrimmed into clips of 20 seconds each and generate descriptions for each clip. These clip-level descriptions are summarized using GPT3.5 Turbo to obtain a video-level description. However, this summarization step loses valuable information and hence fails to provide an accurate summary of the long video. Future work should explore an effective sampling strategy for long video understanding. Another limitation of LLAVIDAL is its diminished efficacy when both Pose and object cues are integrated within the LLVM framework. We believe that modality progressive training could potentially address this suboptimal performance which will be explored in future work.

## J Licensing and Intended Use

This paper introduces a large-scale dataset, **ADL-X**, comprising 100K untrimmed RGB video-instruction pairs, 3D poses, language descriptions, and action-conditioned object trajectories. The raw videos in ADL-X comprise content from NTURGB+D [49], for which the original authors retain distribution rights for the clipped action videos. The scripts utilized to curate the dataset are open-sourced, facilitating the regeneration of the dataset. We will also provide comprehensive features, including image features extracted using CLIP, pose features derived from PoseCLIP, and object features obtained through ObjectLM. We plan to release ADL-X via an academic website for research, academic, and commercial use. The dataset is protected under the **CC-BY** license of Creative Commons, which allows users to distribute, remix, adapt, and build upon the material in any medium or format, as long as the creator is attributed. The license allows ADL-X for commercial use. As the authors of this manuscript and collectors of this dataset, we reserve the right to distribute the data. Additionally, we provide the code, data, and instructions needed to reproduce the main experimental baseline results, and the statistics pertinent to the dataset. We specify all the training details (e.g., data splits, hyperparameters, model-specific implementation details, compute resources used, etc.). Furthermore, we release the code and model weights of our proposed **Large LAnguage VIision** model for **Daily Activities of Living (LLAVIDAL)**, along with the features and instruction QA pairs for the combination videos. The ADL-X dataset focuses on ADL and does not contain any personal data that can resemble evidence, reveal identification, or show offensive content.

The ADL-X dataset can be used by multiple domain experts to advance research and development in various applications related to ADL. Its potential applications include, but are not limited to, assistive technologies, healthcare monitoring systems [73], smart homes [74], robotics for assisted living, and instructional videos for ADL training and support. The dataset can also contribute to the development of AI-driven solutions that aim to improve the quality of life for individuals with disabilities, older adults, and those in need of daily assistance. While we believe that the ADL-X dataset has the potential to make a positive impact on society by enabling the development of technologies that support and enhance the lives of individuals, we acknowledge that, as with any technology, there is a possibility that the dataset or the ideas it presents could be misused or adapted for harmful purposes. However, as authors, we strongly oppose any detrimental usage of this dataset, regardless of whether it is by an individual or an organization, under profit or non-profit motivations. We pledge not to support any endeavors that could cause harm to individuals or society in relation to our data or the ideas presented herein. Our intention is to foster research and innovation in the field of ADL analysis and support, ultimately contributing to the development of technologies that improve the quality of life for those who need assistance with daily activities. We encourage all users of the ADL-X dataset to adhere to the highest ethical standards and to prioritize the well-being of individuals and society in their research and development efforts.