# HourVideo

## 1-Hour Video-Language Understanding

*NeurIPS 2024 Datasets and Benchmarks*

hourvideo.stanford.edu

**Keshigeyan Chandrasegaran**, Agrim Gupta, Lea M. Hadzic, Taran Kota, Jimming He,
Cristobal Eyzaguirre, Zane Durante, Manling Li, Jiajun Wu, Li Fei-Fei

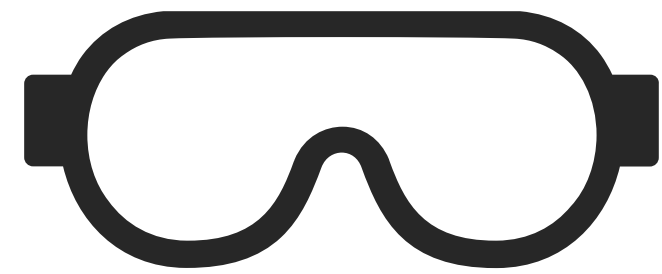Long-form Video-Language Understanding

How can I get to the storage room from the living room?

Where did I leave my AirPods after working out?

Lack of holistic methods for evaluating hour-long video understanding capabilities.

Summarization

Perception

Visual Reasoning

Navigation

500 egocentric
videos
(381 hours)

77 everyday
scenarios

18 total
tasks

12,976
high quality
MCQs

# [HourVideo Design] Proposed Task Suite

## Summarization
- Key Events/ Objects Identification
- Temporal Sequencing
- Compare / Contrast

## Perception

| Information Retrieval | Factual Recall |
| | Sequence Recall |
| | Temporal Distance |
| | Tracking |

## Navigation
- Room-to-Room
- Object Retrieval

## Visual Reasoning

### Spatial
- Relationship
- Proximity
- Layout

### Temporal
- Duration
- Frequency
- Pre-requisites

### Predictive

### Causal

### Counterfactual

## Summarization

Key Events/ Objects Identification

Temporal Sequencing

Compare / Contrast

## Perception

Information Retrieval

Factual Recall

Sequence Recall

Temporal Distance

Tracking

## Navigation

Room-to-Room

Object Retrieval

We manually design question prototypes for each task/ sub-task.

## Visual Reasoning

### Spatial

Relationship

Proximity

Layout

### Temporal

Duration

Frequency

Pre-requisites

### Predictive

### Causal

### Counterfactual

**1** Video Curation



> 100 hours

**1** Video Curation

> 100 hours

**2** MCQ Generation

LLM

$MCQ_2$

**①** Video Curation



> 100 hours

**②** MCQ Generation

LLM

$MCQ_2$

**③** MCQ Refinement using Human Feedback

> 400 hours

LLM

$MCQ_2$

$MCQ_3$

**1** Video Curation

> 100 hours

**2** MCQ Generation

LLM

$MCQ_2$

**3** MCQ Refinement using Human Feedback

$MCQ_2$

> 400 hours

LLM

$MCQ_3$

**4** Blind Filtering

$MCQ_3$

LLM

$MCQ_4$

**1** Video Curation

> 100 hours

**2** MCQ Generation

LLM

$MCQ_2$

**3** MCQ Refinement using Human Feedback

$MCQ_2$

> 400 hours

LLM

$MCQ_3$

**4** Blind Filtering

$MCQ_3$

LLM

$MCQ_4$

**5** Expert MCQ Refinement

$MCQ_4$

> 300 hours

$MCQ_5$

# HourVideo Statistics



**1**

- Cooking
- Cleaning / laundry
- Construction / renovation
- Eating
- Crafting/knitting
- Carpenter
- Talking with family members
- Indoor Navigation (walking)
- Watching tv
- On a screen (phone/laptop)
- Listening to music
- Grocery shopping indoors
- Playing with pets
- Walking on street
- Gardening
- Baker
- Doing yardwork
- Car - commuting, road trip
- Bike mechanic
- Working at desk

(x-axis: 0, 70, 140)

**2**

| Summarization (714) | |
|---|---|
| Key Events/ Objects Identification | 467 |
| Temporal Sequencing | 152 |
| Compare/Contrast | 95 |

| Perception (3777) | |
|---|---|
| Factual Recall | 2479 |
| Sequence Recall | 854 |
| Temporal Distance | 267 |
| Tracking | 177 |

| Navigation (312) | |
|---|---|
| Room-to-Room | 120 |
| Object Retrieval | 192 |

| Spatial (3173) | |
|---|---|
| Relationship | 1889 |
| Proximity | 1239 |
| Layout | 45 |

| Temporal (4292) | |
|---|---|
| Duration | 1945 |
| Frequency | 1815 |
| Pre-requisites | 532 |

Predictive (407)

Causal (150)

Counterfactual (151)

**3** Duration (in minutes)

(Count histogram: 20-30, 30-40, 40-50, 50-60, 60-70, 70-80, 80-90, 90-100, 100-110, 110-120)

**4** #MCQs per video

(Count histogram)

## Summarization

Key Events/ Objects Identification

Temporal Sequencing

Compare / Contrast

## Perception

| Information Retrieval | Factual Recall |
| | Sequence Recall |
| | Temporal Distance |
| Tracking | |

## Navigation

Room-to-Room

Object Retrieval

## Visual Reasoning

| Spatial | Temporal | Predictive |
| --- | --- | --- |
| Relationship | Duration | |
| Proximity | Frequency | Causal |
| Layout | Pre-requisites | Counterfactual |

Identify the unique individuals the camera wearer interacted with.

00:00:00

00:09:13

00:03:55

00:18:23

00:04:04

00:30:00

MCQ Test

☐ 2 Adults

☐ 1 Adult

☐ 4 Adults

☐ 5 Adults

☑ 3 Adults
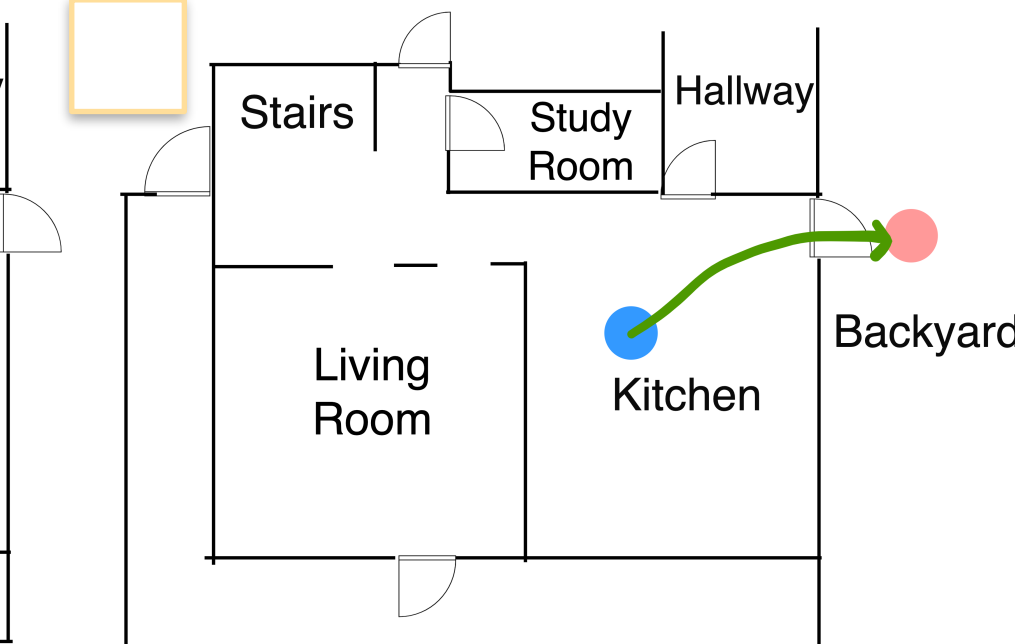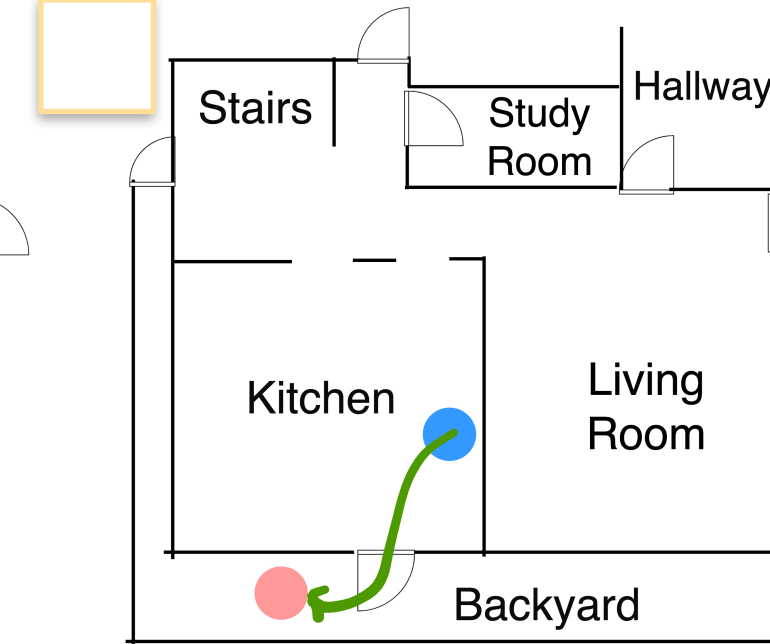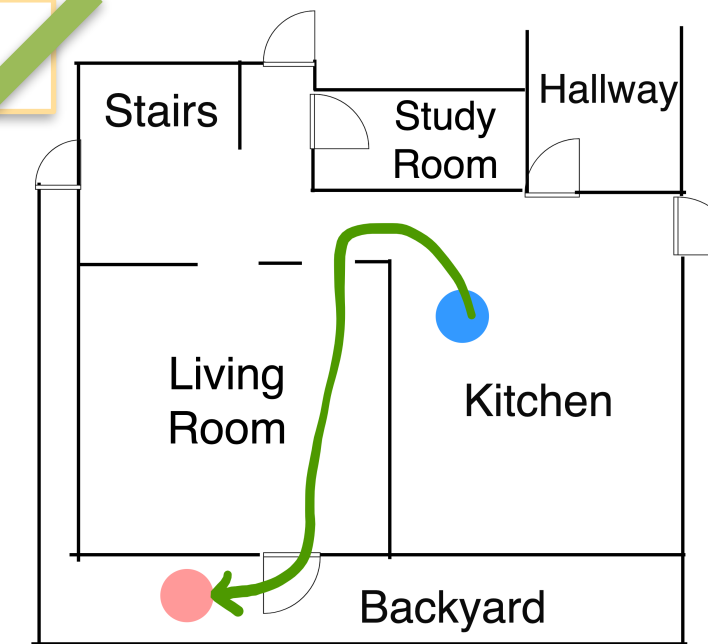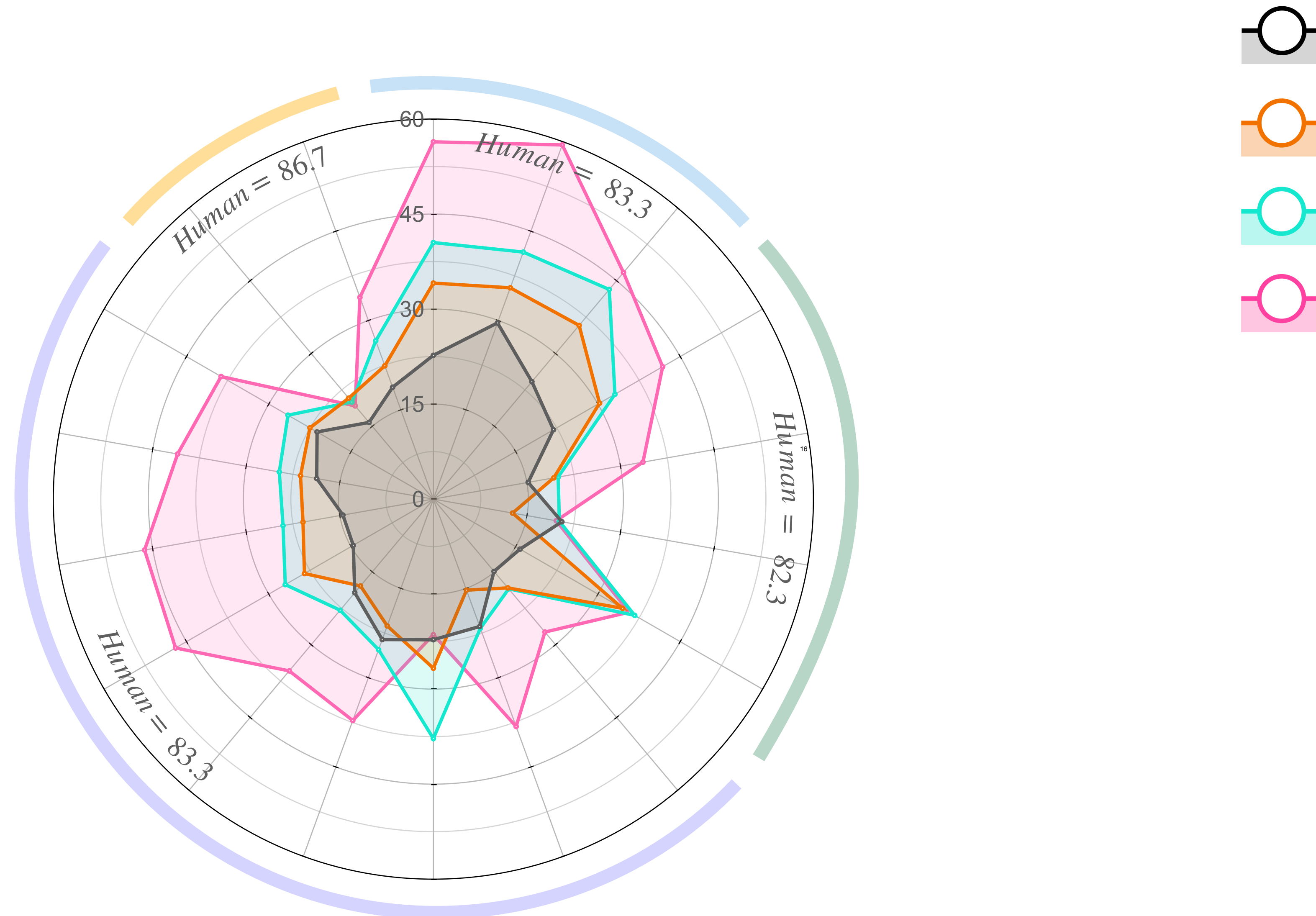
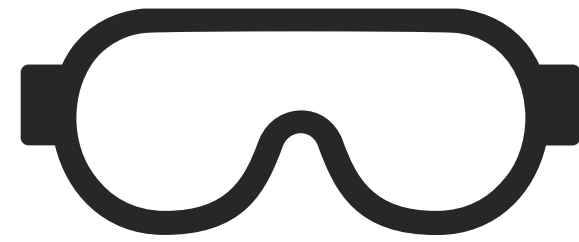How can the camera wearer get to the backyard from the kitchen?

MCQ Test

How do today's multimodal models perform on the HourVideo Benchmark?

# Key Takeaways

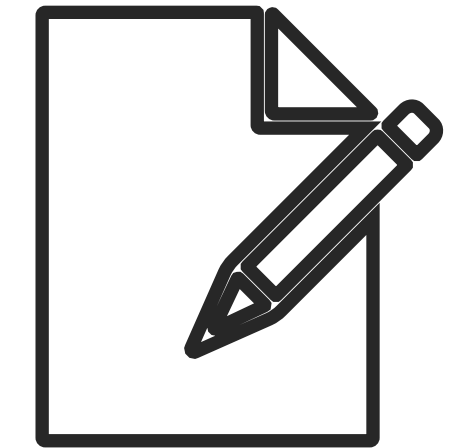500 egocentric
videos
(381 hours)

77 everyday
scenarios

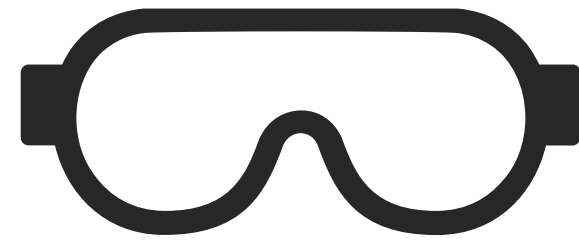Summarization

Perception

Visual Reasoning

Navigation

18 total
tasks

12,976
high quality
MCQs

Summarization

Perception

Visual Reasoning

Navigation

500 egocentric
videos
(381 hours)

77 everyday
scenarios

18 total
tasks

12,976
high quality
MCQs

We introduce **HourVideo**, a benchmark dataset designed to rigorously evaluate the capabilities of multimodal models to comprehend hour-long videos.

Summarization

Perception

Visual Reasoning

Navigation

500 egocentric
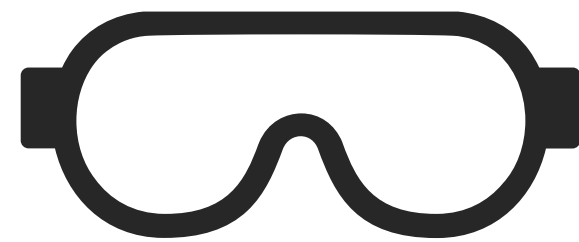videos
(381 hours)

77 everyday
scenarios

18 total
tasks

12,976
high quality
MCQs

We introduce **HourVideo**, a benchmark dataset designed to rigorously evaluate the capabilities of multimodal models to comprehend hour-long videos.

We show that **a significant gap** exists between human experts and SOTA multimodal foundation models in comprehending long-form videos.

Paper/ Benchmark/ Code/ Demos

hourvideo.stanford.edu