

# Action100M: A Large-scale Video Action Dataset

Delong Chen<sup>1,2</sup>, Tejaswi Kasarla<sup>1,3</sup>, Yejin Bang<sup>1</sup>, Mustafa Shukor<sup>1,4</sup>, Willy Chung<sup>1,4</sup>, Jade Yu<sup>1</sup>, Allen Bolourchi<sup>1</sup>, Théo Moutakanni<sup>1</sup>, Pascale Fung<sup>1,2</sup>

<sup>1</sup>Meta FAIR, <sup>2</sup>HKUST, <sup>3</sup>University of Amsterdam, <sup>4</sup>Sorbonne Université

Inferring physical actions from visual observations is a fundamental capability for advancing machine intelligence in the physical world. Achieving this requires large-scale, open-vocabulary video action datasets that span broad domains. We introduce ACTION100M, a large-scale dataset constructed from 1.2M Internet instructional videos (14.6 years of duration), yielding  $O(100)$  million temporally localized segments with open-vocabulary action supervision and rich captions. ACTION100M is generated by a fully automated pipeline that (i) performs hierarchical temporal segmentation using V-JEPA 2 embeddings, (ii) produces multi-level frame and segment captions organized as a TREE-OF-CAPTIONS, and (iii) aggregates evidence with a reasoning model (GPT-OSS-120B) under a multi-round SELF-REFINE procedure to output structured annotations (brief/detailed action, actor, brief/detailed caption). Training VL-JEPA on ACTION100M demonstrates consistent data-scaling improvements and strong zero-shot performance across diverse action recognition benchmarks, establishing Action100M as a new foundation for scalable research in video understanding and world modeling.

Correspondence: [delong.chen@connect.ust.hk](mailto:delong.chen@connect.ust.hk), [theomoutakanni@meta.com](mailto:theomoutakanni@meta.com)

Dataset: <https://github.com/facebookresearch/Action100M>



## 1 Introduction

Making machine intelligence useful in the physical world requires AI models that not only understand world states (*e.g.*, objects and their attributes), but also recognize **physical actions** that interact with the world and induce state transitions. Powered by supervision from large datasets (Laurençon et al., 2023; Awadalla et al., 2024; Shukor et al., 2025; Schuhmann et al., 2022), world state perception in frontier models has advanced rapidly (Alayrac et al., 2022; Chen et al., 2022; Bai et al., 2023; Liu et al., 2024; Dai et al., 2023; Marafioti et al., 2025). In contrast, the capability of understanding *actions* is comparably limited, largely due to the absence of robust data foundations. Existing video action datasets are mainly separately developed for individual narrow domains (*e.g.*, cooking, toy assembly), and remain in limited scale (*e.g.*, less than 1 million action instances).

Advancement in data foundation will enable the development of **open-domain** and **open-vocabulary** video action recognizers. These models will support embodied learning, facilitate wearable assistive applications, and advance physical world modeling (Chen et al., 2025b,a; Fung et al., 2025; Terver et al., 2025; Ball et al., 2025; Staff, 2024=5; Russell et al., 2025; Agarwal et al., 2025a; Xiang et al., 2024). Achieving this necessitates training data that is sufficiently *large scale*, maintains *high quality*, and spans *broad diversity*.

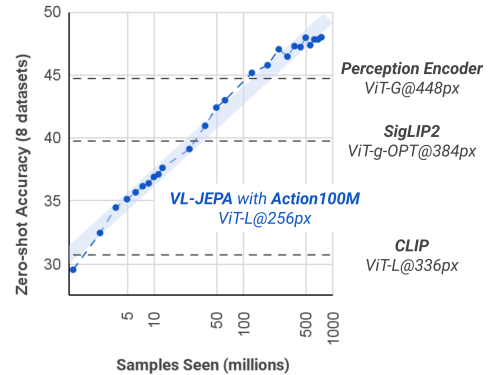
To address these challenges, we introduce ACTION100M, a large-scale video action dataset containing  $O(100)$  million action instances, annotated via a fully automated pipeline. The construction process begins with online instructional videos (Miech et al., 2019), which capture people interacting with the physical world in diverse activities. Dense action labels are produced using a family of frontier open-source models, including V-JEPA 2 (Assran et al., 2025), PerceptionLM (Cho et al., 2025), Llama-3.2-Vision, and GPT-OSS (Agarwal et al., 2025b). The TREE-OF-CAPTIONS (Chen et al., 2024a) and SELF-REFINE mechanisms (Chen et al., 2025b) are leveraged to reduce hallucinations and ensure annotation quality. This pipeline yields a hierarchy of temporal segments annotated with structured fields, including brief/detailed action descriptions, the actor, and brief/detailed video captions.

The annotation pipeline was executed on 1.2 million YouTube videos (14.6 years of duration), spending approximately 1.3 million V100 GPU hours for segmentation and captioning, and 0.3 million H100/H200 GPU

hours for LLM aggregation. It derives 147 million segment-level annotations, totaling 21.3 billion English words.

To demonstrate the utility of ACTION100M, we train VL-JEPA models (Chen et al., 2025c) on the dataset. VL-JEPA is a vision-language model that extracts visual embeddings with V-JEPA 2 (Assran et al., 2025) encoder and predicts target embedding generated by a text encoder. The model uses InfoNCE loss to learn an aligned embedding space that allows CLIP-style open-vocabulary classification (Radford et al., 2021). Evaluation is conducted on eight downstream video action benchmarks, including Something-something-v2 (Goyal et al., 2017), EPIC-KITCHENS-100 (Damen et al., 2020), EgoExo4D Keysteps (Grauman et al., 2024), Kinetics-400 (Kay et al., 2017), COIN (Tang et al., 2019), and CrossTask (Zhukov et al., 2019). These benchmarks evaluate actions spanning diverse domains and different levels of abstraction, from low-level fine-grained ones like “turn on light”, “cut tomato”, to high-level procedural ones like “make Irish coffee”, “assemble computer”, etc.

Fig. 1 demonstrates the **effectiveness of scaling** for zero-shot action recognition: VL-JEPA performance improves consistently as the amount of ACTION100M training data increases (x-axis shows effective batch size  $\times$  number of iterations). Our VL-JEPA (ViT-L at 8 frames per input in 256px) outperform CLIP (Radford et al., 2021), SigLIP2 (Tschannen et al., 2025), and Perception Encoder (Bolya et al., 2025) which have seen much more samples (13B, 40B, and 86B, respectively, while VL-JEPA only 3B) and use larger backbone and higher input resolution. At the same time, the caption annotations in ACTION100M allows VL-JEPA to establish robust vision-language alignment that allows effective zero-shot text-to-video retrieval. Compared to CLIP, SigLIP2, and Perception Encoder, VL-JEPA achieves higher average retrieval recall@1 on eight benchmarks: MSR-VTT (Xu et al., 2016), ActivityNet (Caba Heilbron et al., 2015), DiDeMo (Anne Hendricks et al., 2017), MSVD (Chen and Dolan, 2011), YouCook2 (Zhou et al., 2018), PVDBench (Bolya et al., 2025), Dream-1K (Wang et al., 2024a), and VDC-1K (Chai et al., 2024).



**Figure 1** Scaling on ACTION100M improves zero-shot action recognition consistently.

## 2 Related Works

We present a comparison of ACTION100M with existing action and caption datasets in Table 1, and discuss the detailed related work below. We organize prior efforts into two main categories: video action datasets, which focus on annotated action segments, and video captioning datasets, which provide descriptive text for video content.

**Video Action Datasets.** Standard video action datasets, such as COIN (Tang et al., 2019), CrossTask (Zhukov et al., 2019), and YouCook2 (Zhou et al., 2018), have played a pivotal role in advancing action detection in untrimmed videos. These datasets typically segment long activities into semantically meaningful steps, each paired with textual descriptions, and are primarily composed of instructional videos covering procedural tasks. The collection methodologies for video datasets vary significantly:

- Participant-recorded datasets (e.g., Breakfast (Kuehne et al., 2014), Assembly101 (Sener et al., 2022)) provide dense, hierarchical annotations in controlled environments with high-quality labels through manual data collection.
- Internet-mined datasets (e.g., CrossTask, COIN) leverage external taxonomies to achieve broader coverage of activities. By mining videos from online sources, these datasets capture a wider range of human actions.
- Egocentric datasets (e.g., EgoProceL (Bansal et al., 2022), Ego4D Goal-Step (Song et al., 2023)) introduce data-driven, hierarchical taxonomies and dense step annotations, particularly in domains such as cooking. These datasets support research on goal inference and long-term temporal modeling by capturing first-person perspectives and fine-grained activity sequences.

**Table 1 Summary of major video caption and action recognition datasets.** This table summarizes key statistics for prominent video datasets, including total duration, number of videos and clips, average text length, and annotation sources. The top table lists large-scale video caption datasets and the bottom table presents action recognition and instructional video datasets, including both manual and automated annotations. ACTION100M is distinguished by its unprecedented scale (100M clips), rich hierarchical annotations (brief and detailed actions and captions), and broad coverage of real-world activities.

**(a) Action Recognition and Instructional Video Datasets**

Dataset	Duration	#Videos	#Clips	Avg Text Length	Annotation
COIN (Tang et al., 2019)	476 hours	11.8K	46.3K	4.8	Manual
YouCook2 (Zhou et al., 2018)	176 hours	2K	14K	8.8	Manual
THUMOS14 (Idrees et al., 2017)	30 hours	2,584	20,108	–	Manual
ActivityNet Captions (Caba Heilbron et al., 2015)	849 hours	20K	100K	13.5	Manual
FineAction (Liu et al., 2022)	705 hours	17K	103K	–	–
EGTEA (Sudhakaran et al., 2021)	28 hours	86	10,325	–	Manual
50Salads (Stein and McKenna, 2013)	4 hours	–	–	–	–
Breakfast (Kuehne et al., 2014)	77 hours	–	11,267	–	Manual
Assembly101 (Sener et al., 2022)	513 hours	4321	1M	–	Manual
EgoProceL (Bansal et al., 2022)	62 hours	329	–	–	Manual
Ego4D-Goal-step (Song et al., 2023)	368 hours	851	48K	–	Manual
Action100M Brief Action	14.6 years	1.2M	147M	18.4	PLM-3B, Llama-3.2-Vision-11B, GPT-OSS-120B
Action100M Detailed Action	14.6 years	1.2M	147M	150.2	PLM-3B, Llama-3.2-Vision-11B, GPT-OSS-120B

**(b) Caption Datasets**

Dataset	Duration	#Videos	#Clips	Avg Text Len	Annotation
YT-Temporal-180M (Zellers et al., 2021)	–	6M	180M	–	ASR
HD-VILA-100M (Xue et al., 2022)	42.4 years	3.3M	103M	32.5	ASR
InternVid (Wang et al., 2024c)	86.8 years	7.1M	234M	17.6	Tag2Text, BLIP2
VidChapters-7M (Yang et al., 2023)	35.1 years	817K	6.8M	–	ASR, user chapters
Panda-70M (Chen et al., 2024c)	19.0 years	70.8M	–	13.2	VideoLlama, MiniGPT-4, etc
ShareGPT4Video (Chen et al., 2024b)	291 hours	40K	–	273.3	GPT-4V
OpenVid-1M (Nan et al., 2025)	0.23 years	1M	–	126.5	LLaVA-v1.6-34B
MiraData (Ju et al., 2024)	1.8 years	330K	–	318	GPT-4V
VidGen-1M (Tan et al., 2024)	0.34 years	1M	–	89.3	VILA, Llama-3.1
Koala-36M (Wang et al., 2025)	19.6 years	36M	–	202.1	GPT-4V
EgoVid-5M (Wang et al., 2024b)	0.63 years	5M	–	–	LLaVA-NeXTVideo, Qwen2
OpenHumanVid (Li et al., 2025)	1.9 years	13.2M	–	–	Llama 3.1, BLIP2
VideoUFO (Wang and Yang, 2025)	0.44 years	568K	1M	155.5	Qwen2-VL-7B
UltraVideo (Xue et al., 2025)	62 hours	5K	42K	824.2	Qwen2.5-VL-72B
PE Video Dataset (Cho et al., 2025)	0.49 years	1M	120K	111.7	PLM, Llama-3.3-70B, human refined
PLM-Video-Auto (Cho et al., 2025)	6.06 years	6.4M	–	–	Llama-3.3-70B, LLaMA-3-405B
Action100M Brief Caption	14.6 years	1.2M	147M	106.8	PLM-3B, Llama-3.2-Vision-11B, GPT-OSS-120B
Action100M Detailed Caption	14.6 years	1.2M	147M	540.0	PLM-3B, Llama-3.2-Vision-11B, GPT-OSS-120B

Despite these advances, existing datasets are often constrained by manual annotation bottlenecks, limited domain coverage, or lack of scale.

While datasets such as COIN (Tang et al., 2019) and YouCook2 (Zhou et al., 2018) provide valuable annotated instructional videos, they are constrained by manual annotation and limited activity coverage. Internet-mined datasets expand diversity but often lack fine-grained, hierarchical labels. Egocentric datasets introduce new perspectives but remain domain-specific. ACTION100M marks a step forward in this landscape.

In contrast, ACTION100M leverages large-scale online instructional videos (Miech et al., 2019) to deliver over 100 million action instances with multi-level, open-vocabulary annotations. This enables robust research in open-domain action recognition and world modeling, as summarized in Table 1. By addressing annotation bottlenecks and supporting fine-grained temporal reasoning, ACTION100M provides the way for future advances in video understanding.

**Video Caption Datasets.** A number of large-scale video–language datasets have been proposed in recent years. InternVid (Wang et al., 2024c), HD-VILA-100M (Xue et al., 2022), and VidChapters-7M (Yang et al., 2023) collect millions of videos or hundreds of millions of clips with associated textual descriptions. However, the majority of these captions are obtained from noisy sources such as automatic speech recognition (ASR)

transcripts or video metadata. As a result, the captions are typically short, generic, and only weakly aligned with the underlying physical actions, limiting their usefulness for learning fine-grained action representations or detailed world models.

To address the limitations of ASR-only captions, recent instruction-style datasets leverage powerful vision-language models to produce richer descriptions. OpenVid-1M (Nan et al., 2025), UltraVideo (Xue et al., 2025), and VideoUFO (Wang and Yang, 2025) generate longer, more detailed, and instruction-like captions that better capture object interactions and step-by-step activities. Nonetheless, these datasets are typically built on a relatively small number of source videos compared to the largest ASR-based corpora, which constrains their coverage of diverse environments and long-horizon activities. Koala-36M (Wang et al., 2025) further scales LLM-captioned video data, but is primarily oriented toward video generation rather than structured action understanding.

For video understanding and perception-centric applications, PE Video and PLM-Video-Auto (Cho et al., 2025) combine millions of videos with LLM-generated captions. While they improve caption quality over pure ASR pipelines, they generally assign a single caption per segment and do not explicitly encode temporal hierarchies within videos. Consequently, they provide limited supervision for modeling multi-scale structure in activities (e.g., steps, sub-tasks, and overarching tasks), which is crucial for learning world models that reason over extended sequences of actions and their context.

Moreover, many existing video captioning datasets rely heavily on manual or semi-manual annotation, which is expensive and restricts scale, or they prioritize scale at the cost of fine-grained, action-centric detail and temporal structure. In contrast, ACTION100M is constructed to combine the scale of ASR-driven corpora with the descriptiveness of recent LLM-captioned datasets. It introduces a Tree-of-Captions structure derived from hierarchical temporal segmentation, providing multiple levels of captions that capture both fine-grained actions and broader contextual narratives. This design preserves temporal hierarchy in a dense-captioning-like manner, making ACTION100M particularly suited for training large-scale world models and action-centric video understanding systems.

### 3 Action100M Data Pipeline

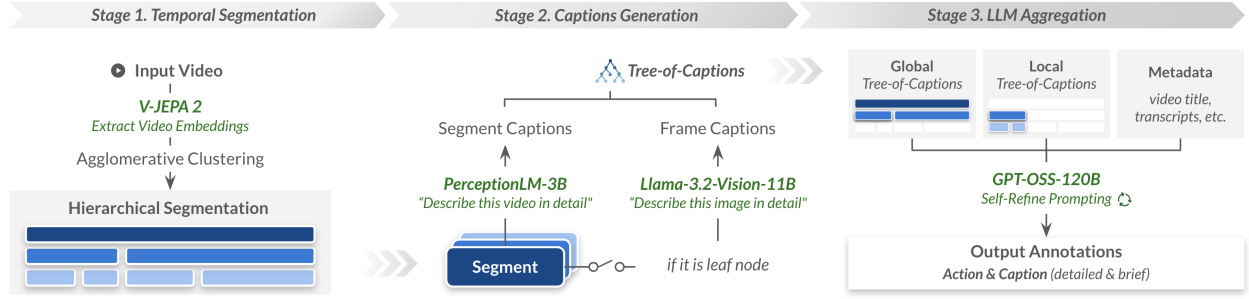
This section describes an automated, scalable pipeline for constructing ACTION100M. The pipeline extends and improves the data generation procedure introduced in Vision Language World Model (Chen et al., 2025b). Fig. 2 provides an overview of the pipeline, and Fig. 3 visualizes an example of annotation. The pipeline has three stages:

1. We decompose each video into a hierarchy of temporally coherent segments, which captures both short, fine-grained motions and long, procedural steps, enabling action supervision at multiple levels of abstraction.
2. For every segment, we generate complementary frame- and video-level captions and organize them into a TREE-OF-CAPTIONS. Captioning is a task for which modern VLMs are extensively trained, and even medium-sized models can produce reliable results.
3. We prompt an LLM to aggregate information from the TREE-OF-CAPTIONS and to extract structured annotations. Assembling evidence across multiple caption levels provably reduces hallucinations (Chen et al., 2024a).

Overall, rather than applying heavy VLMs directly to entire videos, our pipeline first converts videos into hierarchical, text-based representations (TREE-OF-CAPTIONS, inspired by Pyramid-of-Captions (Chen et al., 2024a)) and only leverages strong reasoning models in the final stage, which is purely text-based. This design allows us to obtain reliable annotations while keeping the overall computation cost manageable. In the following, we describe each stage in detail.

**Stage 1. Temporal Segmentation.** Each video is first transformed into a temporally dense sequence of visual embeddings using the V-JEPA 2 encoder. Frames are uniformly sampled at one out of every four raw frames to approximate the temporal resolution used during V-JEPA 2 pretraining. The video is then divided into overlapping temporal windows, each containing 64 sampled frames, with an eight-frame stride between





**Figure 2 Action100M Data Pipeline.** Our pipeline first applies hierarchical temporal segmentation to decompose the video into semantically coherent segments at multiple temporal scales. For each segment, we generate video caption and frame captions, capturing both temporal and spatial information. Next, we prompt LLM to aggregate the captions, extracting final annotations.

consecutive windows. Each window is independently processed by the V-JEPA 2 ViT-g-384<sup>1</sup> encoder, which outputs a sequence of spatial-temporal visual tokens. We perform spatial average pooling, resulting in a per-frame feature of dimensionality equal to the encoder’s hidden size. Because adjacent windows overlap, multiple representations are produced for shared frames. These are accumulated and averaged to form a single, temporally consistent embedding per frame across the entire video.

To capture actions across multiple temporal scales, we apply *hierarchical agglomerative clustering* to the sequence of frame-level representations<sup>2</sup>. A local temporal connectivity constraint links each frame only to its immediate neighbors, ensuring that merges occur only between contiguous time spans. Clustering proceeds bottom-up using Ward linkage, which minimizes the variance within each segment at every merge step. Starting from individual frames or short windows, adjacent segments are iteratively merged to minimize intra-cluster variance. This process produces a hierarchical tree decomposition of the video, where each node corresponds to a contiguous, semantically coherent temporal segment. Lower levels of the hierarchy correspond to fine-grained atomic motions, while higher levels capture coarser activities. We retain only nodes whose duration is larger than 0.5 seconds, ensuring that each segment is semantically meaningful.

**Stage 2. Caption Generation.** After hierarchical segmentation, we annotate each node in the video tree with captions, ensuring that both local fine-grained and global contextual information are captured within the same representation. The captioning process operates in two complementary modes: mid-frame captioning for fine-grained spatial details and video-segment captioning for temporal dynamics.

For every leaf node (representing the smallest contiguous action segment), we extract a key frame at the midpoint of its temporal span. These mid-frames are processed using Llama-3.2-Vision-11B<sup>3</sup>, prompted with “Describe this image in detail.” to generate frame captions. For higher-level nodes representing longer temporal spans, we apply Perception-LM-3B<sup>4</sup>. Each segment is sampled into 32 evenly spaced frames between its start and end times at 320<sup>2</sup> resolution, and the model is prompted with “Describe this video in detail.” to generate segment-level captions. Both models are configured with a generation limit of 1024 tokens and can be run on a single NVIDIA V100 32GB GPU.

**Stage 3. LLM Aggregation.** For each node in the TREE-OF-CAPTIONS, we generate structured action annotations by prompting GPT-OSS-120B<sup>5</sup> to extract five fields of information: brief action description, detailed action description, actor, brief video caption, and detailed video caption. This extraction is performed through three iterative rounds of SELF-REFINE to improve consistency and quality.

The inputs to the LLM include the current node’s caption, its children’s captions formatted in depth-first order, as well as global context such as root captions (within a limited depth) and video metadata (including title, description, and ASR transcript). Each node in the TREE-OF-CAPTIONS is processed independently. Nodes

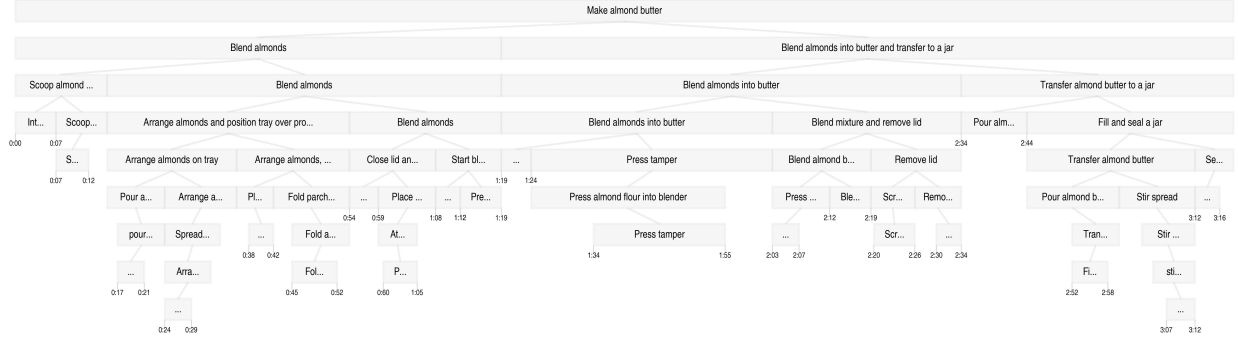
<sup>1</sup><https://huggingface.co/facebook/vjepa2-vitg-fpc64-384>

<sup>2</sup><https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html>

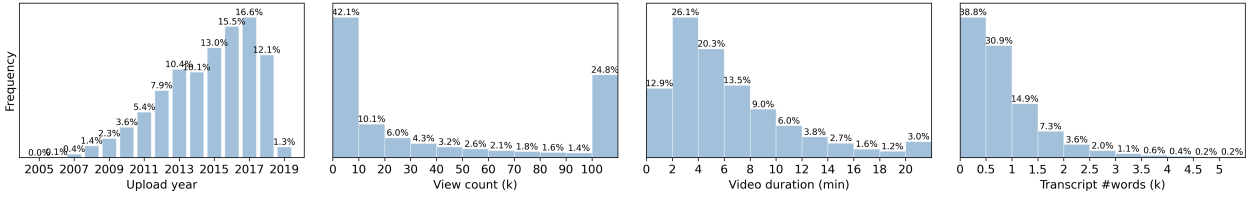
<sup>3</sup><https://huggingface.co/meta-llama/Llama-3.2-11B-Vision>

<sup>4</sup><https://huggingface.co/facebook/Perception-LM-3B>

<sup>5</sup><https://huggingface.co/openai/gpt-oss-120b>



**Figure 3** Example of hierarchical structure in Action100M annotations (with brief action description labels shown). Source video: [url](#). **Brief caption** of the entire video: *A woman roasts almonds, blends them into butter, and pours the butter into a jar.* **Detailed caption:** *The video opens with the presenter in a bright kitchen speaking to the camera. She spreads raw almonds on a parchment-lined tray, places the tray in a pre-heated 350 °F oven, and after roasting lets the nuts cool to room temperature. She then transfers the almonds to a Vitamix blender, removes the lid, inserts a tamper, and sets the machine on high. While the blender runs she presses the almonds down with the tamper, first creating a fine flour and then a thick creamy butter within about one minute. She pours the almond butter into a clear storage jar, scoops it with a large wooden spoon and stirs it to smooth the surface, then concludes the segment with a brief thank-you.*



**Figure 4** Statistics of Action100M source videos and metadata. Distributions of (left to right) video upload year, view count, video duration, and transcript length, computed over the subset of videos for which metadata is available.

with a duration shorter than four seconds are discarded. For the remaining nodes, we query GPT-OSS-120B to infer clean, structured textual representations that unify and denoise information from multiple caption sources. More details about the process and the prompt to the GPT-OSS-120B are in the supplementary.

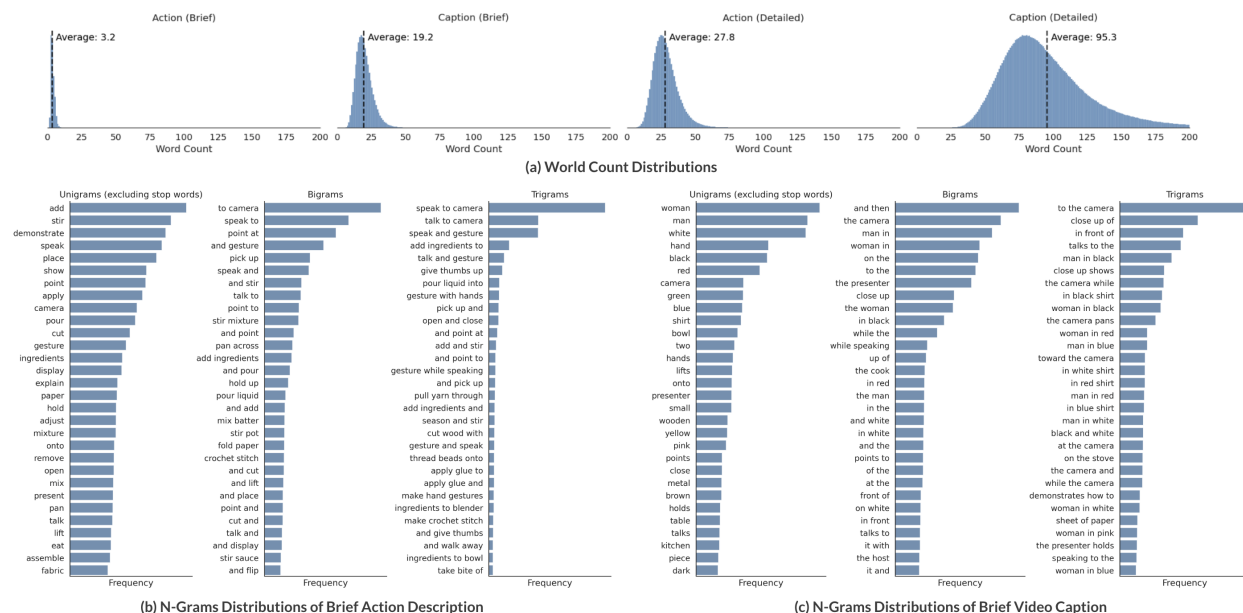
## 4 Dataset Analysis

**Source Videos.** Action100M is built on 1,199,096 face-blurred videos sourced from HowTo100M (Miech et al., 2019), corresponding to a total video duration of approximately 14.6 years. Since the original HowTo100M dataset was released in June 2019 (Miech et al., 2019), many videos have become unavailable since then. We successfully retrieved ASR transcripts for 72% of these videos, covering 10.6 years of video content. Figure 4 summarizes the statistics of the available metadata. Figure 5 visualizes a word cloud derived from video titles. As HowTo100M focuses on instructional content curated from 12 WikiHow categories (e.g., Food & Entertaining, Home & Garden, Hobbies & Crafts, while excluding more abstract categories such as Relationships or Finance), the most frequent keywords strongly reflect procedural and hands-on activities. Common terms include “make”, “recipe”, “DIY”, “easy”, “cake”, and “chocolate”, highlighting the dominance of cooking, home improvement, and everyday skill-oriented videos in the dataset.

**Generated Annotations.** Figure 6(a) summarizes the word-count statistics of the four annotation types produced by our pipeline. The average length increases monotonically from brief action (3.2 words), to brief caption (19.2 words), to detailed action (27.8 words), and finally to detailed caption (95.3 words). Across the entire dataset, Action100M contains 147,092,653 annotated video segments (including videos without metadata), corresponding to an estimated total of 0.46B (brief actions) + 2.83B (brief captions) + 3.96B (detailed actions) + 14.02B (detailed captions) = 21.27B words.



**Figure 5 Word cloud of video titles in Action100M.** Frequently occurring words reflect the instructional and procedural nature of the dataset, with dominant terms related to cooking, DIY activities, and everyday physical tasks.



**Figure 6 Statistics of generated textual annotations in Action100M.** (a) Word-count distributions for four annotation types: brief action, brief caption, detailed action, and detailed caption. Dashed lines indicate the mean length of each annotation type. (b-c) Top unigrams, bigrams, and trigrams (excluding stop words) in brief action descriptions and brief video captions.

Due to the hierarchical temporal segmentation procedure, shorter segments are substantially more common than longer ones. Specifically, 64% of all segments have durations between 0-3 seconds, followed by 23.8% in the 3-10 second range, and 10.2% between 10 seconds and 1 minute. Only about 2% of segments are longer than one minute. For action annotations, 3.23% of segments are labeled as “N/A” by GPT-OSS-120B, which typically correspond to non-action content such as video introductions, advertisements, or subscription reminders. Storing all annotations together with metadata and the full Tree-of-Captions structure requires approximately 205 GB of disk space.

Figures 6(b) and (c) show the N-gram distributions of brief action descriptions and brief video captions, respectively. As expected, action descriptions are dominated by verbs (e.g., *add*, *stir*, *demonstrate*), whereas video captions contain a higher proportion of adjectives and object-centric descriptors. The frequency distributions further reveal a strong imbalance in action concepts, with certain patterns such as “*speak to camera*” occurring disproportionately often. This observation motivates the semantic resampling strategy introduced in §5.4 to alleviate long-tail imbalance during training. Figure 7 provides a qualitative view of this structure through sunburst diagrams for the five most frequent words in brief action descriptions.



**Figure 7 Sunburst visualizations of frequent action compositions.** Each sunburst shows the hierarchical co-occurrence structure centered on one of the five most frequent verbs in brief action descriptions.

## 5 Experiments

### 5.1 Training and Evaluation Setup

We train VL-JEPA model (Chen et al., 2025c). We perform the query-free pretraining (*i.e.*, BASE) with ACTION100M data. We refer to (Chen et al., 2025c) for details of model architecture. We train VL-JEPA in three stages (Tab. 2):

- In **Stage 1**, we perform image pretraining with single frame per input. We use DataComp-1B (Gadre et al., 2023), mixing original caption and model generated caption from Li et al. (2024), and YFCC-100M (Thomee et al., 2016) captioned by PaliGemma (Beyer et al., 2024) and Florence-2 (Xiao et al., 2024).
- In **Stage 2**, we process video data with eight frames per input. We continue from stage 1 checkpoint and use ACTION100M data (mixing all four field, with detailed action and caption downsampled by a half), and PerceptionLM-3B labelled action description and video caption for each ACTION100M segment.
- In **Stage 3**, we increase the number of frames per input to 32, and unfreeze the V-JEPA 2 encoder. This leads to higher CUDA memory consumption and lower batch size, and we use gradient accumulation of 4 steps and lower learning rate (1e-5 instead of 5e-5).

**Table 2** Details of VL-JEPA’s three-stage training procedure. Stage 3 uses a gradient accumulation of 4 steps.

Training Stage	Vision encoder	#Frames per input	Training Data	Batch Size	#Iterations
Stage 1	Frozen	1	Image-text data (DataComp-1B, YFCC-100M)	24,576	100k
Stage 2	Frozen	8	Action100M	12,288	60k
Stage 3	Unfrozen	32	Action100M	3,072×4	10k

We benchmark resulting models with existing foundation models, including CLIP (Radford et al., 2021), SigLIP2 (Tschannen et al., 2025), and Perception Encoder (Bolya et al., 2025). We evaluate on two tasks:

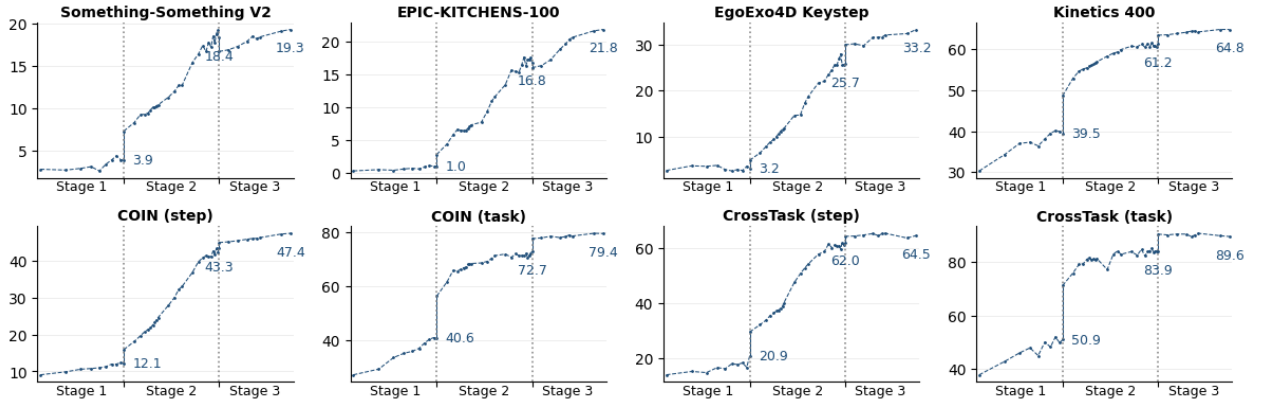
- **Zero-shot action recognition** (top-1 accuracy) on eight benchmarks: Something-something-v2 (SSv2) (Goyal et al., 2017), EPIC-KITCHENS-100 (EK-100) (Damen et al., 2020), EgoExo4D Keysteps (Grauman et al., 2024), Kinetics-400 (Kay et al., 2017), COIN (Tang et al., 2019), and CrossTask (Zhukov et al., 2019). For COIN and Crosstask, we evaluate both segment-level step recognition and global task recognition.
- **Zero-shot text-to-video retrieval** (recall@1) on eight benchmarks: MSR-VTT (Xu et al., 2016), ActivityNet (Caba Heilbron et al., 2015), DiDeMo (Anne Hendricks et al., 2017), MSVD (Chen and Dolan, 2011), YouCook2 (Zhou et al., 2018), PVD-Bench (Bolya et al., 2025), Dream-1K (Wang et al., 2024a), and VDC-1K (Chai et al., 2024).

### 5.2 Main Results

As shown in Tab. 3, with ACTION100M, VL-JEPA achieves higher average action recognition performance and average video retrieval performance. Thanks to strong V-JEPA 2 encoder and dense fine-grained action

**Table 3 Main results.** We evaluate zero-shot performance on eight action recognition dataset. and eight video retrieval datasets.

Model		#Parameters	#Samples Seen	#Frames	Action Recognition (Top-1 Accuracy)								Text-to-video Retrieval (Recall@1)									
					Average	SSv2	EK100	EgoExo4D	Kinetics-400	COIN (SR)	COIN (TR)	CrossTask (SR)	CrossTask (TR)	Average	MSR-VTT	ActivityNet	DiDeMo	MSVD	YouCook2	PVD-Bench	Dream-1k	VDC-1k
CLIP	RN50 (224px)	75M			21.8	2.1	1.5	2.1	41.4	8.6	39.0	10.9	68.7	28.3	28.7	17.7	24.7	29.7	5.1	27.6	47.2	46.0
	ViT-B (224px)	124M	12.8B	8	25.4	3.1	1.3	2.8	49.5	11.2	47.3	16.2	71.5	29.3	31.0	19.5	25.7	34.0	6.1	27.0	48.5	42.9
	ViT-L (336px)	389M			30.7	3.8	3.7	2.6	58.3	14.7	63.5	20.8	78.5	35.3	35.9	23.4	30.7	41.9	7.9	36.7	56.8	49.3
SigLIP2	ViT-B (224px)	375M			33.9	5.2	2.3	4.5	57.8	20.6	69.9	27.7	82.9	39.6	40.2	25.0	32.1	48.6	13.8	52.1	60.9	43.7
	ViT-L (384px)	882M	40B	8	38.6	5.9	4.5	6.4	63.6	24.2	78.5	35.1	90.8	45.4	41.6	32.7	35.1	53.5	19.0	59.2	71.6	50.9
	ViT-g (384px)	1.9B			39.8	6.1	6.1	5.6	68.0	26.0	80.4	35.1	90.8	47.5	43.4	33.9	38.9	56.0	22.2	60.4	73.0	52.5
PE-Core	ViT-B (224px)	448M	58B		37.2	5.8	3.3	6.0	65.4	21.5	77.1	26.9	91.8	44.9	46.5	35.4	35.3	49.1	15.2	59.8	68.7	49.2
	ViT-L (336px)	671M	58B	8	42.9	9.3	6.0	11.6	73.4	27.1	83.3	37.5	95.3	50.2	48.9	41.7	40.8	56.2	22.5	64.7	75.9	51.0
	ViT-G (448px)	2.3B	86B		44.7	9.0	6.4	13.6	<b>76.4</b>	29.0	<b>86.0</b>	40.3	<b>97.2</b>	58.1	<b>51.6</b>	49.1	44.5	<b>58.7</b>	26.0	77.0	89.2	68.5
VL-JEPA	Stage 1		2.4B	8	21.5	3.9	1.0	3.2	39.5	12.1	40.6	20.9	50.9	32.6	27.7	20.9	26.6	34.2	3.4	51.6	52.9	43.8
	Stage 2		3.1B	8	48.0	18.4	16.8	25.7	61.2	43.3	72.7	62.0	83.9	59.5	37.1	54.2	46.0	45.7	33.9	80.2	91.5	87.3
	Stage 3		3.3B	32	<b>52.5</b>	<b>19.3</b>	<b>21.8</b>	<b>33.2</b>	64.8	<b>47.4</b>	79.4	<b>64.5</b>	89.6	<b>63.7</b>	40.0	<b>64.9</b>	<b>50.0</b>	49.0	<b>40.4</b>	<b>83.1</b>	<b>93.3</b>	<b>88.8</b>



**Figure 8 Zero-shot action recognition accuracy in each training stage.** Stage 1: image pretraining with single frame input; Stage 2: main Action100M pretraining with 8 frames input; Stage 3: 32 frames input with unfrozen encoder. The x-axis of each stage represents number of training samples seen (*i.e.*, number of iterations) in log-scale. The performance at the end of each stage is annotated.

annotation in ACTION100M, VL-JEPA is particularly strong on **motion-focused** tasks, such as Something-something-v2, EPIC-KITCHENS-100, EgoExo4D Keysteps, step recognition on COIN and CrossTask.

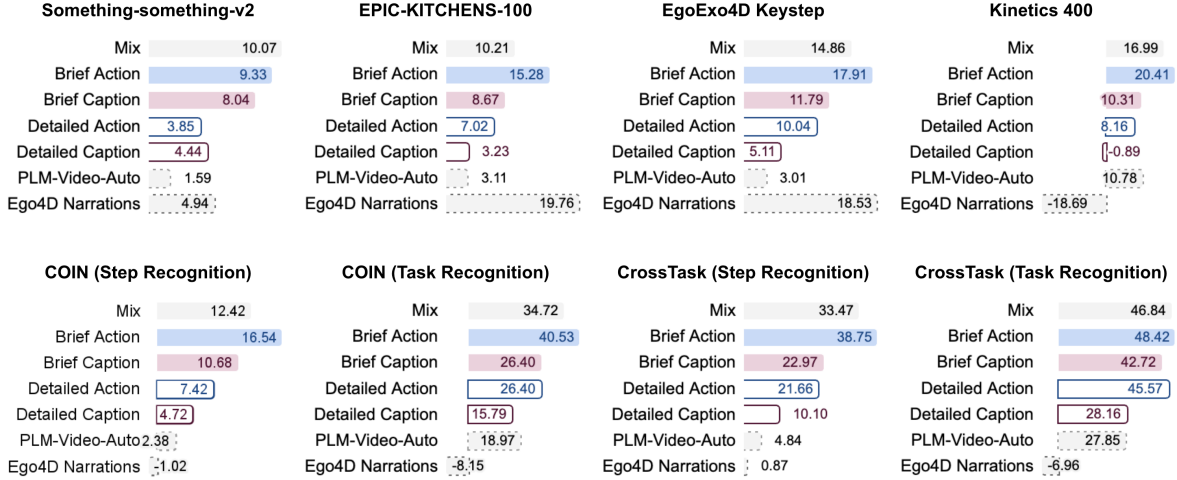
Fig. 8 visualize evolution of per-dataset accuracy in each stages, with x-axis being log-scale number of samples seen in each stage. We see that scaling yields consistent improvement, especially on motion-focused datasets mentioned earlier. We also see a significant jump from stage 1 to stage 2, indicating that image-only training alone is insufficient for action recognition.

### 5.3 Effectiveness of Action100M Pipeline

We compare the effectiveness of different data under a controlled training setting. All models are initialized from the same stage 1 checkpoint, and further trained for 20k steps (20.48M samples seen). Fig. 9 visualizes the resulting performance improvement compared to the stage 1 initialization.

We see that brief action descriptions in ACTION100M are highly effective in terms of improving zero-shot action recognition. It outperform the direct PLM-3B pseudo labeling baseline, showing the effectiveness of hierarchical captioning and LLM aggregation. Detailed captions in ACTION100M also show advantages over PLM-Video-Auto captions on most benchmarks. At the same time, ACTION100M is much larger than it (100M vs 3M). Training with Ego4D atomic action description significantly improve egocentric action recognition performance on EK-100 and EgoExo4D, while not being effective on other domains.





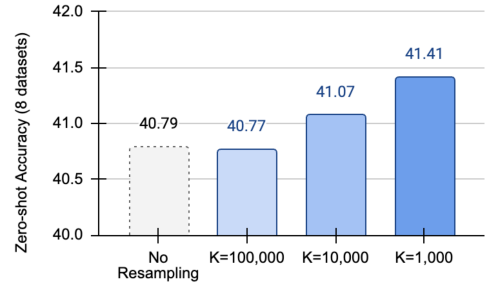
**Figure 9** Performance improvements of stage 2 video pretraining with different data upon stage 1.

## 5.4 Effectiveness of Semantic Resampling

A major challenge in large-scale action datasets is the long-tailed distribution of action frequencies, which can bias model training. To address this, we explore semantic resampling inspired by DINOv2 (Vo et al., 2024) to promote a more balanced sampling of actions.

We begin by embedding all brief action descriptions from the TREE-OF-CAPTIONS using EmbeddingGemma-300M (Vera et al., 2025). To address redundancy, we deduplicate these actions by hashing their text, ensuring that repeated descriptions are consolidated. Following deduplication, we apply k-means clustering to the resulting embeddings, with cluster counts of  $k = \{10^3, 10^4, 10^5\}$ . This clustering groups semantically similar actions together, allowing us to control the granularity of the action space. From each cluster, actions are sampled uniformly with replacement until the target dataset size is reached, which ensures that both frequent and rare actions are adequately represented in the training data.

To empirically validate the effectiveness of this pipeline, we curated a 10M-action dataset using semantic resampling across all tested values of  $k$ . All models were initialized from the same stage 1 checkpoint and subsequently trained for 10k steps, corresponding to 10.24M samples seen. We report the average across the 8 benchmark zero-shot action classification datasets. As shown in Fig. 10, resampling with a smaller number of clusters leads to improved performance compared to training without resampling, highlighting the benefits of down-sampling frequent actions and up-sampling rare ones.



**Figure 10** Effectiveness of semantic resampling.

## 6 Conclusion

We introduced ACTION100M, a large-scale, open-domain dataset for action-centric video understanding built from 1.2M procedural videos and annotated into 147M temporally localized segments. Empirically, pretraining VL-JEPA on ACTION100M yields consistent scaling behavior and strong zero-shot transfer across a diverse set of action recognition and text-to-video retrieval benchmarks, with notable strengths on motion-focused and step-centric datasets. Controlled ablations further show the value of LLM-aggregated brief actions and hierarchical evidence over direct pseudo-labeling, and semantic resampling improves sample efficiency by mitigating long-tail redundancy. Overall, ACTION100M provides a practical route to scale open-vocabulary action understanding, and it enables future work on action anticipation, action-conditioned world models, and long-horizon planning that requires reasoning over multi-scale procedural structure.

## References

- Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025a.
- Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv preprint arXiv:2508.10925*, 2025b.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *Proceedings of the IEEE international conference on computer vision*, pages 5803–5812, 2017.
- Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zhohus, et al. V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *arXiv preprint arXiv:2506.09985*, 2025.
- Anas Awadalla, Le Xue, Oscar Lo, Manli Shu, Hannah Lee, Etash Guha, Sheng Shen, Mohamed Awadalla, Silvio Savarese, Caiming Xiong, et al. Mint-1t: Scaling open-source multimodal data by 10x: A multimodal dataset with one trillion tokens. *Advances in Neural Information Processing Systems*, 37:36805–36828, 2024.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. 2023. <https://api.semanticscholar.org/CorpusID:261101015>.
- Philip J. Ball, Jakob Bauer, Frank Belletti, Bethanie Brownfield, Ariel Ephrat, Shlomi Fruchter, Agrim Gupta, Kristian Holsheimer, Aleksander Holynski, Jiri Hron, Christos Kaplanis, Marjorie Limont, Matt McGill, Yanko Oliveira, Jack Parker-Holder, Frank Perbet, Guy Scully, Jeremy Shar, Stephen Spencer, Omer Tov, Ruben Villegas, Emma Wang, Jessica Yung, Cip Baetu, Jordi Berbel, David Bridson, Jake Bruce, Gavin Buttmore, Sarah Chakera, Bilva Chandra, Paul Collins, Alex Cullum, Bogdan Damoc, Vibha Dasagi, Maxime Gazeau, Charles Gbadamosi, Woohyun Han, Ed Hirst, Ashyana Kachra, Lucie Kerley, Kristian Kjems, Eva Knoepfel, Vika Koriakin, Jessica Lo, Cong Lu, Zeb Mehring, Alex Moufarek, Henna Nandwani, Valeria Oliveira, Fabio Pardo, Jane Park, Andrew Pierson, Ben Poole, Helen Ran, Tim Salimans, Manuel Sanchez, Igor Saprykin, Amy Shen, Sailesh Sidhwani, Duncan Smith, Joe Stanton, Hamish Tomlinson, Dimple Vijaykumar, Luyu Wang, Piers Wingfield, Nat Wong, Keyang Xu, Christopher Yew, Nick Young, Vadim Zubov, Douglas Eck, Dumitru Erhan, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Raia Hadsell, Aaron van den Oord, Inbar Mosseri, Adrian Bolton, Satinder Singh, and Tim Rocktäschel. Genie 3: A new frontier for world models. 2025.
- Siddhant Bansal, Chetan Arora, and CV Jawahar. My view is the best view: Procedure learning from egocentric videos. In *European Conference on Computer Vision*, pages 657–675. Springer, 2022.
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun Cho, Andrea Madotto, Chen Wei, Tengyu Ma, Jiale Zhi, Jathushan Rajasegaran, Hanoona Rasheed, et al. Perception encoder: The best visual embeddings are not at the output of the network. *arXiv preprint arXiv:2504.13181*, 2025.
- 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.
- Wenhao Chai, Enxin Song, Yilun Du, Chenlin Meng, Vashisht Madhavan, Omer Bar-Tal, Jenq-Neng Hwang, Saining Xie, and Christopher D Manning. Auroracap: Efficient, performant video detailed captioning and a new benchmark. *arXiv preprint arXiv:2410.03051*, 2024.
- David Chen and William B Dolan. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies*, pages 190–200, 2011.
- Delong Chen, Samuel Cahyawijaya, Etsuko Ishii, Ho Shu Chan, Yejin Bang, and Pascale Fung. What makes for good image captions? *arXiv preprint arXiv:2405.00485*, 2024a.

- Delong Chen, Willy Chung, Yejin Bang, Ziwei Ji, and Pascale Fung. Worldprediction: A benchmark for high-level world modeling and long-horizon procedural planning. *arXiv preprint arXiv:2506.04363*, 2025a.
- Delong Chen, Theo Moutakanni, Willy Chung, Yejin Bang, Ziwei Ji, Allen Bolourchi, and Pascale Fung. Planning with reasoning using vision language world model. *arXiv preprint arXiv:2509.02722*, 2025b.
- Delong Chen, Mustafa Shukor, Theo Moutakanni, Willy Chung, Jade Yu, Tejaswi Kasarla, Allen Bolourchi, Yann LeCun, and Pascale Fung. Vl-jepa: Joint embedding predictive architecture for vision-language. *arXiv preprint arXiv:2512.10942*, 2025c.
- Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Zhenyu Tang, Li Yuan, et al. Sharegpt4video: Improving video understanding and generation with better captions. *Advances in Neural Information Processing Systems*, 37:19472–19495, 2024b.
- Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13320–13331, 2024c.
- Xi Chen, Xiao Wang, Soravit Changpinyo, Anthony J 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.
- Jang Hyun Cho, Andrea Madotto, Effrosyni Mavroudi, Triantafyllos Afouras, Tushar Nagarajan, Muhammad Maaz, Yale Song, Tengyu Ma, Shuming Hu, Suyog Jain, et al. Perceptionlm: Open-access data and models for detailed visual understanding. *arXiv preprint arXiv:2504.13180*, 2025.
- 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. *Advances in neural information processing systems*, 36:49250–49267, 2023.
- Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. The epic-kitchens dataset: Collection, challenges and baselines. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11):4125–4141, 2020.
- Pascale Fung, Yoram Bachrach, Asli Celikyilmaz, Kamalika Chaudhuri, Delong Chen, Willy Chung, Emmanuel Dupoux, Hongyu Gong, Hervé Jégou, Alessandro Lazaric, et al. Embodied ai agents: Modeling the world. *arXiv preprint arXiv:2506.22355*, 2025.
- Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36:27092–27112, 2023.
- Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The "something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pages 5842–5850, 2017.
- Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19383–19400, 2024.
- Haroon Idrees, Amir R Zamir, Yu-Gang Jiang, Alex Gorban, Ivan Laptev, Rahul Sukthankar, and Mubarak Shah. The thumos challenge on action recognition for videos “in the wild”. *Computer Vision and Image Understanding*, 155:1–23, 2017.
- Xuan Ju, Yiming Gao, Zhaoyang Zhang, Ziyang Yuan, Xintao Wang, Ailing Zeng, Yu Xiong, Qiang Xu, and Ying Shan. Miradata: A large-scale video dataset with long durations and structured captions. *Advances in Neural Information Processing Systems*, 37:48955–48970, 2024.
- Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- Hilde Kuehne, Ali Arslan, and Thomas Serre. The Language of Actions: Recovering the Syntax and Semantics of Goal-Directed Human Activities . In *2014 IEEE Conference on Computer Vision and Pattern Recognition*

- (CVPR), pages 780–787, Los Alamitos, CA, USA, June 2014. IEEE Computer Society. doi: 10.1109/CVPR.2014.105. <https://doi.ieeecomputersociety.org/10.1109/CVPR.2014.105>.
- Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, et al. Obelics: An open web-scale filtered dataset of interleaved image-text documents. *Advances in Neural Information Processing Systems*, 36:71683–71702, 2023.
- Hui Li, Mingwang Xu, Yun Zhan, Shan Mu, Jiaye Li, Kaihui Cheng, Yuxuan Chen, Tan Chen, Mao Ye, Jingdong Wang, et al. Openhumanvid: A large-scale high-quality dataset for enhancing human-centric video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 7752–7762, 2025.
- Xianhang Li, Haoqin Tu, Mude Hui, Zeyu Wang, Bingchen Zhao, Junfei Xiao, Sucheng Ren, Jieru Mei, Qing Liu, Huangjie Zheng, et al. What if we recaption billions of web images with llama-3? *arXiv preprint arXiv:2406.08478*, 2024.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26296–26306, 2024.
- Yi Liu, Limin Wang, Yali Wang, Xiao Ma, and Yu Qiao. Fineaction: A fine-grained video dataset for temporal action localization. *IEEE transactions on image processing*, 31:6937–6950, 2022.
- Andrés Marafioti, Orr Zohar, Miquel Farré, Merve Noyan, Elie Bakouch, Pedro Cuenca, Cyril Zakka, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, et al. Smolvlm: Redefining small and efficient multimodal models. *arXiv preprint arXiv:2504.05299*, 2025.
- 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 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2630–2640, 2019.
- Kepan Nan, Rui Xie, Penghao Zhou, Tiehan Fan, Zhenheng Yang, Zhijie Chen, Xiang Li, Jian Yang, and Ying Tai. Openvid-1m: A large-scale high-quality dataset for text-to-video generation. In *The Thirteenth International Conference on Learning Representations*, 2025. <https://openreview.net/forum?id=j7kdXSrISM>.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.
- Lloyd Russell, Anthony Hu, Lorenzo Bertoni, George Fedoseev, Jamie Shotton, Elahe Arani, and Gianluca Corrado. Gaia-2: A controllable multi-view generative world model for autonomous driving. *arXiv preprint arXiv:2503.20523*, 2025.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022.
- Fadime Sener, Dibyadip Chatterjee, Daniel Shelepov, Kun He, Dipika Singhania, Robert Wang, and Angela Yao. Assembly101: A large-scale multi-view video dataset for understanding procedural activities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21096–21106, 2022.
- Mustafa Shukor, Louis Bethune, Dan Busbridge, David Grangier, Enrico Fini, Alaaeldin El-Nouby, and Pierre Ablin. Scaling laws for optimal data mixtures. *arXiv preprint arXiv:2507.09404*, 2025.
- Yale Song, Eugene Byrne, Tushar Nagarajan, Huiyu Wang, Miguel Martin, and Lorenzo Torresani. Ego4d goal-step: Toward hierarchical understanding of procedural activities. *Advances in Neural Information Processing Systems*, 36: 38863–38886, 2023.
- World Labs Technical Staff. Generating worlds. <https://www.worldlabs.ai/blog>, 2024=5.
- Sebastian Stein and Stephen J. McKenna. Combining embedded accelerometers with computer vision for recognizing food preparation activities. In *Proceedings of the 2013 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, UbiComp ’13, page 729–738, New York, NY, USA, 2013. Association for Computing Machinery. ISBN 9781450317702. doi: 10.1145/2493432.2493482. <https://doi.org/10.1145/2493432.2493482>.
- Swathikiran Sudhakaran, Sergio Escalera, and Oswald Lanz. Learning to recognize actions on objects in egocentric video with attention dictionaries. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):6674–6687, 2021.

- Zhiyu Tan, Xiaomeng Yang, Luo Zheng Qin, and Hao Li. Vidgen-1m: A large-scale dataset for text-to-video generation. *arXiv preprint arXiv:2408.02629*, 2024.
- Yansong Tang, Dajun Ding, Yongming Rao, Yu Zheng, Danyang Zhang, Lili Zhao, Jiwen Lu, and Jie Zhou. Coin: A large-scale dataset for comprehensive instructional video analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1207–1216, 2019.
- Basile Terver, Tsung-Yen Yang, Jean Ponce, Adrien Bardes, and Yann LeCun. What drives success in physical planning with joint-embedding predictive world models? *arXiv preprint arXiv:2512.24497*, 2025.
- Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016.
- Michael Tschanen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*, 2025.
- Henrique Schechter Vera, Sahil Dua, Biao Zhang, Daniel Salz, Ryan Mullins, Sindhu Raghuram Panyam, Sara Smoot, Iftexhar Naim, Joe Zou, Feiyang Chen, et al. Embeddinggemma: Powerful and lightweight text representations. *arXiv preprint arXiv:2509.20354*, 2025.
- Huy V Vo, Vasil Khalidov, Timothée Darcet, Théo Moutakanni, Nikita Smetanin, Marc Szafraniec, Hugo Touvron, Camille Couprie, Maxime Oquab, Armand Joulin, et al. Automatic data curation for self-supervised learning: A clustering-based approach. *arXiv preprint arXiv:2405.15613*, 2024.
- Jiawei Wang, Liping Yuan, Yuchen Zhang, and Haomiao Sun. Tarsier: Recipes for training and evaluating large video description models. *arXiv preprint arXiv:2407.00634*, 2024a.
- Qiheng Wang, Yukai Shi, Jiarong Ou, Rui Chen, Ke Lin, Jiahao Wang, Boyuan Jiang, Haotian Yang, Mingwu Zheng, Xin Tao, et al. Koala-36m: A large-scale video dataset improving consistency between fine-grained conditions and video content. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8428–8437, 2025.
- Wenhao Wang and Yi Yang. VideoUFO: A million-scale user-focused dataset for text-to-video generation. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2025. <https://openreview.net/forum?id=wwlwRuKle7>.
- Xiaofeng Wang, Kang Zhao, Feng Liu, Jiayu Wang, Guosheng Zhao, Xiaoyi Bao, Zheng Zhu, Yingya Zhang, and Xingang Wang. Egovid-5m: A large-scale video-action dataset for egocentric video generation. *arXiv preprint arXiv:2411.08380*, 2024b.
- Yi Wang, Yanan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Yu Qiao. Internvid: A large-scale video-text dataset for multimodal understanding and generation. In *The Twelfth International Conference on Learning Representations*, 2024c. <https://openreview.net/forum?id=MLBdiWu4Fw>.
- Jiannan Xiang, Guangyi Liu, Yi Gu, Qiyue Gao, Yuting Ning, Yuheng Zha, Zeyu Feng, Tianhua Tao, Shibo Hao, Yemin Shi, et al. Pandora: Towards general world model with natural language actions and video states. *arXiv preprint arXiv:2406.09455*, 2024.
- Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4818–4829, 2024.
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016.
- Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu, and Baining Guo. Advancing high-resolution video-language representation with large-scale video transcriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5036–5045, 2022.
- Zhucun Xue, Jiangning Zhang, Teng Hu, Haoyang He, Yanan Chen, Yuxuan Cai, Yabiao Wang, Chengjie Wang, Yong Liu, Xiangtai Li, et al. Ultravideo: High-quality uhd video dataset with comprehensive captions. *arXiv preprint arXiv:2506.13691*, 2025.
- Antoine Yang, Arsha Nagrani, Ivan Laptev, Josef Sivic, and Cordelia Schmid. Vidchapters-7m: Video chapters at scale. *Advances in Neural Information Processing Systems*, 36:49428–49444, 2023.



- Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. Merlot: Multimodal neural script knowledge models. *Advances in neural information processing systems*, 34:23634–23651, 2021.
- Luowei Zhou, Chenliang Xu, and Jason Corso. Towards automatic learning of procedures from web instructional videos. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Dimitri Zhukov, Jean-Baptiste Alayrac, Ramazan Gokberk Cinbis, David Fouhey, Ivan Laptev, and Josef Sivic. Cross-task weakly supervised learning from instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3537–3545, 2019.

## A Implementation Details

We process each node in the TREE-OF-CAPTIONS independently to generate structured action annotations. Nodes shorter than four seconds are discarded. For the remaining nodes, we query a large reasoning model (GPT-OSS-120B<sup>6</sup>) to infer clean, structured textual representations that unify and denoise information from multiple caption sources.

For each node, we construct an input prompt that combines: formatted TREE-OF-CAPTIONS of the current node and the global root node, and video metadata context-including title, description, ASR transcript. TREE-OF-CAPTIONS are formatted as a depth-first traversal into a Markdown-style text stream. The serialized input is then processed by GPT-OSS-120B under a 3 rounds of SELF-REFINE procedure. In the first round, the model performs structured reasoning with a **high** reasoning effort setting to produce an initial draft of the action summary. In subsequent rounds, the same model revisits the previous output together with the original context, iteratively correcting factual errors, resolving inconsistencies, and removing unsupported statements. The prompt to the GPT-OSS-120B is as follows:

```
# Video metadata
```

```
{video_metadata}
```

```
# Global video context
```

```
{formatted_global_tree_of_captions}
```

```
# Current segment to be processed
```

```
{formatted_current_tree_of_captions}
```

```
# Your Task
```

```
The task is to extract structured information from a video segment. The segment spans from {start_time} to {end_time} seconds, within a larger video that runs from {global_start_time} to {global_end_time} seconds. You will be provided with hierarchical captions generated by models operating on local windows or frames. These captions may contain errors or hallucinations. Your goal is to carefully aggregate the information from the captions to produce an accurate, coherent description of this specific segment.
```

Guidelines and requirements:

- Focus on what is visually observable, emphasizing both 1) physical motion, procedural actions, and 2) appearance information, background or text if possible.
- The timestamps in video metadata and markdown titles are generally reliable, but those inside the captions are not.
- Use global captions and video metadata only for disambiguation. Do not add visually unobservable information (e.g., content of speech, names that cannot be inferred from the local video segment) to the results.
- The result covers the current segment ({start\_time} to {end\_time}) only, not the entire video.
- Ignore captions from very short edges at the start or end of the segment if they are semantically discontinuous (e.g., due to scene transitions). Focus on the main central portion of the segment.
- Due to the limited perception of local models and possible hallucinations, there may be inconsistencies among captions.
- Be cautious and conservative, and rely on the majority consensus.
- Think hard. Perform in-depth reasoning to discuss the provided captions and global context, then output a JSON containing final results in plain English text.
- For all fields, use coherent full sentences with proper capitalization and punctuation. Use concise noun phrases without unnecessary qualifiers or parenthetical clarifications.

```
### Task 1. Summarization
```

---

<sup>6</sup><https://huggingface.co/openai/gpt-oss-120b>

Generate both a short informative caption and a comprehensive, dense summary of the video segment.  
Describe events in chronological order, but do not mention any exact timestamps in the summary. Include

### ### Task 2. Action Identification

Identify the main actor and the physical action performed in the current segment. Provide both a brief description that represents the overall action step, and a detailed description that contains sufficient procedural detail. Use "N/A" (without further explanation) if there are no visible actors or physical actions (e.g., static).

#### # Response Formats

## output

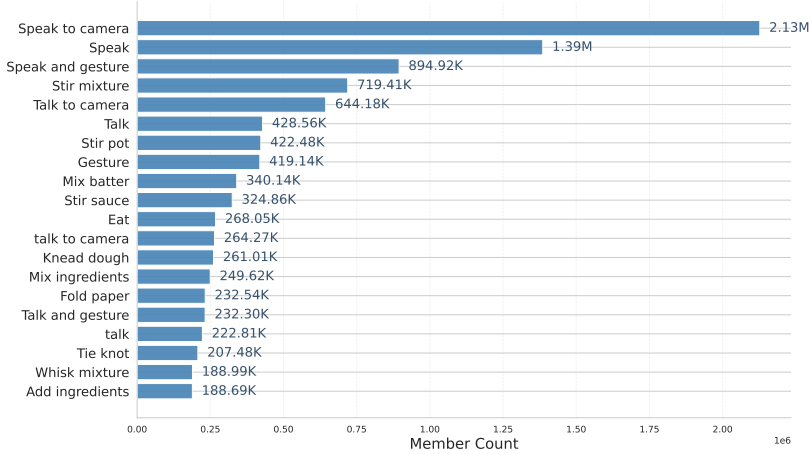
```
{
  "type": "object", "properties": {
    "summary": {"type": "object", "properties": {
      "brief": {
        "type": "string",
        "description": "Single sentence video caption."
      },
      "detailed": {
        "type": "string",
        "description": "Detailed, comprehensive description."
      },
    },
  },
  "action": {"type": "object", "properties": {
    "brief": {
      "type": "string",
      "description": "A single verb phrase (no -ing forms) briefly summarizing the overall action content."
    },
    "detailed": {
      "type": "string",
      "description": "A single imperative sentence describing how the action is performed with more details."
    },
    "actor": {
      "type": "string",
      "description": "Single sentence or an informative noun phrase describing who is performing the action."
    },
  },
},
  "required": ["summary", "action"]
}
```

For SELF-REFINE, we append the following instruction in addition:

Now, carefully analyze, verify, and revise the previous draft so that it is fully accurate, faithful to the provided content, and strictly adheres to all stated guidelines and requirements.

## B Statistics of Duplications and Semantic Resampling

The scale of redundancy in the raw data is substantial. A significant portion of brief action descriptions are repeated, reflecting the inherent long-tailed and redundant nature of large-scale video action data. During the de-duplication step, we identify 7.58 million duplicate groups, which together account for 141.8 million duplicate instances. The remaining action texts are unique, each occurring only once in the dataset. The distribution of these exact duplicates is visualized in Figure 11.



**Figure 11** Distribution of duplicated brief action descriptions.

After deduplication, we apply k-means clustering to the action texts, grouping them into  $k = \{10^3, 10^4, 10^5\}$  clusters. Each cluster represents a set of semantically similar actions, enabling us to control the granularity of the action space. Figure 12 illustrates the effectiveness of this clustering: for each  $k$ , we present examples of selected anchor clusters and their closest neighbors (ranked by cosine similarity), with five random action texts shown per cluster. Across all settings, the clustering pipeline consistently groups together actions with similar semantics, as evidenced by the high cosine similarity between neighboring clusters.

To assess the coverage and diversity of our clustered action space, we analyze the overlap between the  $k = 10^4$  clusters and several standard downstream action recognition datasets. Figure 13 presents a UMAP visualization, where each panel highlights the clusters containing samples from a specific benchmark (COIN, CrossTask, Epic-Kitchens-100, Kinetics-400, YouCook2, EgoExo4D). The varying distributions of colored points across the panels demonstrate that the ACTION100M clusters provide broad and diverse coverage of the action space, with different downstream datasets mapping to distinct but overlapping regions. This analysis confirms that our semantic clustering approach not only reduces redundancy and enhances diversity, but also ensures strong representational alignment with a wide range of downstream tasks.

<b>Anchor • Cluster 184</b> <ul style="list-style-type: none"> <li>• Raise sandwich toward mouth</li> <li>• Eat sandwich</li> <li>• Eat sandwich and sip drink</li> <li>• Serve sandwiches</li> <li>• eat sandwich</li> </ul>	<b>Cluster 276</b> <span>sim 0.850</span> <ul style="list-style-type: none"> <li>• Eat burgers</li> <li>• Flip burgers</li> <li>• Eat burger</li> <li>• Slice cheeseburger</li> <li>• Speed-eat burgers</li> </ul>	<b>Cluster 699</b> <span>sim 0.815</span> <ul style="list-style-type: none"> <li>• Assemble burger</li> <li>• Assemble wedge salad</li> <li>• Assemble salad</li> <li>• Assemble pastry sandwich</li> <li>• Assemble taco salad</li> </ul>	<b>Cluster 141</b> <span>sim 0.751</span> <ul style="list-style-type: none"> <li>• Simmer stew</li> <li>• Smoke pork belly</li> <li>• Sear steak</li> <li>• Pan-broil steak</li> <li>• roast pig</li> </ul>	<b>Cluster 748</b> <span>sim 0.635</span> <ul style="list-style-type: none"> <li>• Light fire</li> <li>• Flail and collapse in a laboratory fire</li> <li>• coat stick with flammable liquid</li> <li>• Cook sausages over fire</li> <li>• Warm hands over fire</li> </ul>
<b>Anchor • Cluster 400</b> <ul style="list-style-type: none"> <li>• Water plants</li> <li>• Place plant</li> <li>• Unwrap plant</li> <li>• pan around plant</li> <li>• Plant runner</li> </ul>	<b>Cluster 462</b> <span>sim 0.908</span> <ul style="list-style-type: none"> <li>• Prune plant</li> <li>• tend garden plants</li> <li>• Water shrub</li> <li>• Mist plants</li> <li>• Weed garden bed</li> </ul>	<b>Cluster 166</b> <span>sim 0.842</span> <ul style="list-style-type: none"> <li>• Examine orchid</li> <li>• Inspect and water sapling</li> <li>• Inspect plants</li> <li>• Inspect and rake garden bed</li> <li>• Inspect plant leaves</li> </ul>	<b>Cluster 852</b> <span>sim 0.751</span> <ul style="list-style-type: none"> <li>• Grab dowel rods</li> <li>• Reach toward tank</li> <li>• Reach toward front fork</li> <li>• Reach into shelf</li> <li>• Grab hat</li> </ul>	<b>Cluster 632</b> <span>sim 0.650</span> <ul style="list-style-type: none"> <li>• Assemble decorative garland</li> <li>• assemble decorative ribbon bows</li> <li>• Assemble flower</li> <li>• cut and glue fabric to assemble a pet bow tie</li> <li>• Assemble beaded heart</li> </ul>

(a) Distribution of texts within clusters for  $k = 10^3$

<b>Anchor • Cluster 134</b> <ul style="list-style-type: none"> <li>• Explain dryer</li> <li>• Explain flash dryer handling</li> <li>• Explain dryer vent function</li> <li>• Explain flashdry</li> <li>• Explain air-drying benefits</li> </ul>	<b>Cluster 5228</b> <span>sim 0.931</span> <ul style="list-style-type: none"> <li>• Explain wood drying</li> <li>• Explain fast drying of burnt umber</li> <li>• Explain grape drying</li> <li>• Explain hardwood floor drying</li> <li>• Explain mushroom drying</li> </ul>	<b>Cluster 56498</b> <span>sim 0.889</span> <ul style="list-style-type: none"> <li>• Explain dryer issue</li> <li>• Explain dryer problem</li> <li>• Explain dryer motor issues</li> <li>• Explain dryer drum assembly</li> <li>• Explain dryer allergen contamination</li> </ul>	<b>Cluster 38735</b> <span>sim 0.785</span> <ul style="list-style-type: none"> <li>• Show hair dryer box</li> <li>• show dryer balls</li> <li>• Show hair dryer</li> <li>• Show dryer balls</li> <li>• Show brush-drying device</li> </ul>	<b>Cluster 55610</b> <span>sim 0.690</span> <ul style="list-style-type: none"> <li>• Explain rice polishing</li> <li>• Explain arsenic in rice</li> <li>• Explain rice consistency</li> <li>• Explain rice sprouting</li> <li>• Explain rice water safety</li> </ul>
<b>Anchor • Cluster 5546</b> <ul style="list-style-type: none"> <li>• Sew pouch</li> <li>• Sew leather pouch</li> <li>• Sew a green pouch</li> <li>• Sew pouch together</li> <li>• Sew fabric pouch</li> </ul>	<b>Cluster 1033</b> <span>sim 0.956</span> <ul style="list-style-type: none"> <li>• Sew pocket onto pouch</li> <li>• Sew green felt piece onto blue pouch</li> <li>• Sew fabric onto pouch</li> <li>• Sew cord into leather pouch</li> <li>• Sew pouch onto cardigan</li> </ul>	<b>Cluster 29113</b> <span>sim 0.917</span> <ul style="list-style-type: none"> <li>• Hand-stitch the orange pouch.</li> <li>• Sew green felt piece onto blue pouch</li> <li>• Hand-stitch pouch</li> <li>• hand-sew pouch</li> <li>• Hand-stitch pocket</li> </ul>	<b>Cluster 68966</b> <span>sim 0.832</span> <ul style="list-style-type: none"> <li>• Pin and sew bag edge</li> <li>• Sew bag edge</li> <li>• Sew border around bag opening</li> <li>• Sew bag edges</li> <li>• Sew bag perimeter</li> </ul>	<b>Cluster 43512</b> <span>sim 0.704</span> <ul style="list-style-type: none"> <li>• Tie stitch</li> <li>• Tie pentagon stitch</li> <li>• Tie knotted stitch</li> <li>• Tie a box stitch</li> <li>• Tie a stitch</li> </ul>

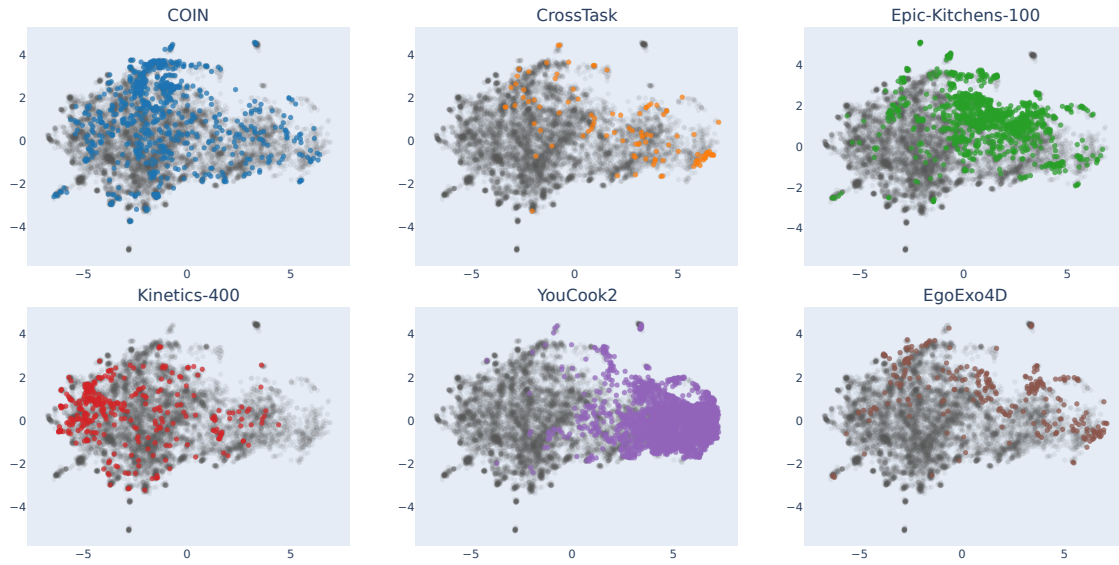
(b) Distribution of texts within clusters for  $k = 10^4$

<b>Anchor • Cluster 15000</b> <ul style="list-style-type: none"> <li>• Place white flower onto gold ring</li> <li>• Push metal ring through petals</li> <li>• Insert flower into ring</li> <li>• Shake flower ring</li> <li>• Place flower on ring</li> </ul>	<b>Cluster 61525</b> <span>sim 0.950</span> <ul style="list-style-type: none"> <li>• Attach flower to ring</li> <li>• attach flower to ring</li> <li>• Attach flowers to wire rings</li> <li>• Attach ring to flower</li> <li>• Wrap tape and attach flower to ring</li> </ul>	<b>Cluster 88225</b> <span>sim 0.858</span> <ul style="list-style-type: none"> <li>• Attach decorative ring with glue</li> <li>• Glue metal ring onto bouquet</li> <li>• Apply glue and attach flowers to the gold ring</li> <li>• Glue ornament onto ring</li> <li>• Glue flower onto ring</li> </ul>	<b>Cluster 99589</b> <span>sim 0.825</span> <ul style="list-style-type: none"> <li>• Place flowers</li> <li>• Place flowers on headboard</li> <li>• place flowers into rain boots</li> <li>• Place red flowers on ice</li> <li>• Place flowers on spheres</li> </ul>	<b>Cluster 39899</b> <span>sim 0.730</span> <ul style="list-style-type: none"> <li>• Attach leather strap to metal ring</li> <li>• Place leather strip and metal ring</li> <li>• Fold leather around ring</li> <li>• Attach leather stopper and metal ring to rope</li> <li>• Thread leather through metal ring</li> </ul>
<b>Anchor • Cluster 87665</b> <ul style="list-style-type: none"> <li>• Wrap tape and pull roll</li> <li>• Wrap cup and pull strip</li> <li>• Wrap strip and pull sock</li> <li>• Wrap and pull cloth</li> <li>• Wrap tape and pull yarn</li> </ul>	<b>Cluster 76472</b> <span>sim 0.918</span> <ul style="list-style-type: none"> <li>• Wrap and pull</li> <li>• Wrap wire and pull thread</li> <li>• Wrap and pull strings</li> <li>• Pull stretch-tie wrap</li> <li>• Wrap strands and pull through hole</li> </ul>	<b>Cluster 80215</b> <span>sim 0.880</span> <ul style="list-style-type: none"> <li>• Pull tape and split roll</li> <li>• Pull and peel duct tape</li> <li>• Pull and cut tape</li> <li>• Pull tape and cut tape</li> <li>• Pull tape and press tool</li> </ul>	<b>Cluster 7438</b> <span>sim 0.835</span> <ul style="list-style-type: none"> <li>• Wrap and cut tape</li> <li>• cut tape and wrap ribbon</li> <li>• Wrap tape and cut it</li> <li>• Cut and wrap tape</li> <li>• Wrap ribbon and cut tape</li> </ul>	<b>Cluster 18971</b> <span>sim 0.755</span> <ul style="list-style-type: none"> <li>• Unwrap food</li> <li>• Unwrap food items</li> <li>• Wrap and heat food</li> <li>• Unwrap food and line tin</li> <li>• Unwrap and arrange wrapped food</li> </ul>

(c) Distribution of texts within clusters for  $k = 10^5$

**Figure 12 Understanding cluster distribution after the semantic clustering (deduplication+kmeans):** Examples of selected anchor clusters (left) and their closest neighboring actions (in decreasing order of similarity, right) for different values of  $k$ , illustrating the effectiveness of the clustering approach. We show 5 random texts within each cluster. *sim* of the neighboring clusters denote the cosine similarity to the anchor cluster. Clusters consistently group together actions with similar semantics across all  $k$  values. Lower  $k$  values ( $10^3$ ) yield broader, more diverse clusters, while higher  $K$  values ( $10^5$ ) produce highly specific clusters.





**Figure 13 UMAP visualization of the overlap between the semantic clusters and downstream datasets.** Each panel containing samples from a specific downstream dataset (colored points) and their overlap with the  $k = 10^4$  cluster of the ACTION100M dataset. This highlights the diversity and coverage of the ACTION100M dataset with respect to multiple downstream benchmarks.