

# MultiVENT 2.0: A Massive Multilingual Benchmark for Event-Centric Video Retrieval

Reno Kriz<sup>♣♥\*</sup> Kate Sanders<sup>♥\*</sup> David Etter<sup>♣♥\*</sup> Kenton Murray<sup>♣♥</sup> Cameron Carpenter<sup>♥</sup>  
 Kelly Van Ochten<sup>♣</sup> Hannah Recknor<sup>♣</sup> Jimena Guallar-Blasco<sup>♣</sup> Alexander Martin<sup>♥</sup>  
 Sydney Johns<sup>♣◇</sup> Ronald Colaianni<sup>♣</sup> Nolan King<sup>♣</sup> Eugene Yang<sup>♣♥</sup> Benjamin Van Durme<sup>♣♥</sup>

<sup>♣</sup> Human Language Technology Center of Excellence <sup>♥</sup> Johns Hopkins University

<sup>♣</sup> SCALE Participants <sup>◇</sup> Virginia Tech University

{rkriz1,ksande25,kenton}@jhu.edu

## Abstract

Efficiently retrieving and synthesizing information from large-scale multimodal collections has become a critical challenge. However, existing video retrieval datasets suffer from scope limitations, primarily focusing on matching descriptive but vague queries with small collections of professionally edited, English-centric videos. To address this gap, we introduce **MultiVENT 2.0**, a large-scale, multilingual event-centric video retrieval benchmark featuring a collection of more than 218,000 news videos and over 3,900 queries targeting specific world events. These queries specifically target information found in the visual content, audio, embedded text, and text metadata of the videos, requiring systems leverage all these sources to succeed at the task. Preliminary results show that state-of-the-art vision-language models struggle significantly with this task, and while alternative approaches show promise, they are still insufficient to adequately address this problem. These findings underscore the need for more robust multimodal retrieval systems, as effective video retrieval is a crucial step towards multimodal content understanding and generation.

## 1. Introduction

While information retrieval systems for text documents have been extensively studied for decades, the landscape has shifted dramatically toward visual content, particularly videos. As of January 2024, YouTube alone likely hosts over 14 billion videos.<sup>1</sup> Despite this explosion of visual data, there remains a dearth of research focused on the efficient retrieval, processing, and synthesis of such vast collections. Datasets that reflect the diverse array of multimodal, multilingual news sources available online could help models adapt to this shift, but existing video retrieval datasets,

<sup>1</sup><https://www.theatlantic.com/technology/archive/2024/01/how-many-videos-youtube-research/677250/>



Figure 1. Example query/video pairs from MSR-VTT [38] and MultiVENT 2.0. MSR-VTT primarily contains broad descriptive queries mapped to general English-centric video clips, while MultiVENT 2.0 targets specific current events covering 4 media formats, 6 languages, and subjects like natural disasters, politics, sports, social gatherings, and science.

including MSR-VTT [38], focus primarily on semantically simple videos created for English-speaking audiences. To address this limitation, the original MultiVENT dataset was introduced, containing **Multilingual Videos of Events** with aligned **Natural Text** across five languages [29]. However, both MultiVENT (2,400 videos) and MSR-VTT (10,000 videos) are extremely small compared to standard information retrieval collections for text documents: In comparison, the HC4 corpus used in the 2022 NEUCLIR TREC shared task contains over 6 million text documents [17]. This disparity in collection size and diversity underscores the challenges in scaling video retrieval research.

Dataset	Query Type	Query Generation	Domain	Source	# Queries	# Videos/(Clips)	Avg Clip Length	Multilingual
ActivityNetCaptions	Caption	Manual	Action	YouTube	100K	20K/100K	120s	no
MSR-VTT	Caption	Manual	Open	YouTube	200K	7.2K/10K	10-30s	no
VaTeX	Caption	Manual	Open/Act	Kinetics-600 (YouTube)	825K	42K/42K	10s	no
DiDeMo	Caption	Manual	Open	YFCC100M (Flickr)	40K	10.5K/26.9K	5s	no
LSMDC	Caption	Manual	Movie	Movies	128K	202/128K	4-5s	no
V3C	Metadata	-	Open	Vimeo	-	28K	8min	yes
MSVD	Caption	Manual	Open	YouTube	70K	2K/2K	4-10s	no
Valor-32K	Caption	Manual	Open	AudioSet (Youtube)	32K	32K/32K	10s	no
YouCook2	Caption	Manual	Cooking	YouTube	15.4K	2K/15.4K	20s	no
TVR	Query	Manual	TV Shows	TV Shows	109K	6/21.8K	76s	no
MTVR	Query	Manual	Movies	TV Shows	218K	6/21.8K	76s	yes
HowTo100M	Subtitle	Manual/Automatic	Instructional	YouTube	136M	1.2M/136M	4s	no
Multi-HowTo100M	Subtitle	Manual/Automatic	Instructional	YouTube	136M	1.2M/136M	9s	yes
MultiVENT 1.0	Caption	Manual	Open	YouTube/Twitter	2,400	2.4K	83s	yes
MULTIVENT 2.0 (ours)	Query	Manual	Open	YouTube/Twitter	3,900	218K	145s	yes

Table 1. Comparison of video retrieval datasets. Many of these datasets were collected as resources for general vision-language model training and did not explicitly target video retrieval. Subsequent use of these collections as a retrieval benchmark leveraged descriptive video captions or subtitles as proxy queries for English only data, which differ greatly in scope and difficulty from cross-lingual text retrieval datasets [17]. To our knowledge, MULTIVENT 2.0 is the first large-scale video dataset with queries explicitly developed for an event-focused multilingual retrieval task.

To push the boundaries of video retrieval and create a more challenging and realistic task, we introduce MULTIVENT 2.0. This dataset consists of more than 218,000 videos and over 3,900 manually-written queries targeting information about specific world events depicted within this video corpus. The videos primarily span six languages—Arabic, Chinese, English, Korean, Russian, and Spanish—and range from professionally edited news broadcasts to raw, first-person footage captured on cell phones. The events in MULTIVENT 2.0 include a wide array of types, including social and sporting events, disasters, political developments, and scientific discoveries. Queries were designed to target specific information from the visual content, audio, embedded text, and text metadata, which challenges models to effectively process and integrate information across modalities and languages and better reflects real-world retrieval scenarios.

Preliminary results show that this task presents significant challenges for current state-of-the-art vision-language models (VLMs). While specialized single-modality models show promise for queries targeting content from their respective modalities, they remain insufficient for addressing the full range of queries. These findings suggest that existing systems are not yet equipped to handle complex vision-language tasks and underscores the need for more robust multimodal systems, as effective video retrieval is a critical step toward multimodal content understanding and generation. The main contributions of this paper are:

1. We introduce MULTIVENT 2.0,<sup>2</sup> a large-scale multilingual video retrieval task containing more than 218,000 videos. Over 3,900 queries were crafted to target aspects of current events using visual content, audio, embedded

text, and text metadata from these videos.

2. We evaluate a variety of video retrieval models, demonstrating that even state-of-the-art multimodal systems struggle in this challenging event-centric setting.
3. We baseline specialized modality-specific models; despite not being able to address the task on their own, we show their potential usefulness as components in more robust retrieval pipelines.

## 2. Related Work

**Video Retrieval Datasets** There are an increasing number of video datasets, many of which support the downstream text-video retrieval task. However, there are various limitations to these existing datasets. Many of these datasets have too few videos or have videos that are too short for comprehensive evaluation of complex event understanding [28]. For example, the popular baseline dataset, MSVD [3] only has 1,970 videos whose length average between 4 and 10 seconds. The larger and more recent VATeX [32] and VALOR-32K datasets [5] have video clips that are on average only 10 seconds long, and the clips of the extremely large dataset of over 1 million videos, HOWTO100M [23], average to only 4 seconds. Many video datasets are limited to videos of a single topic or domain, such as the LSMDC dataset [26] and TVR dataset [18], which feature clips from movies and TV shows respectively, and the YOUCOOK2 dataset which features cooking videos [42]. Thus, while these datasets are quite large in terms of number of videos, they do not represent diverse events or “real-world” topics. Moreover, while some datasets do have non-English videos, they are limited to only one non-English language; MULTIVENT 1.0 [14] and MTVR [19] both have only English and Chinese videos. Furthermore, despite the increase

<sup>2</sup>Data available [at this link](#).

in Video datasets, very few video datasets can support event-retrieval or true multi-modal retrieval. The DIDEMO dataset [12] is a large dataset featuring event-centric videos. However, the creators of the DIDEMO dataset filtered out any non-professional edited videos, thus limiting the scope of the dataset. Additionally, captions or video descriptions are often extracted for pretraining as the “query” intended for the text-video retrieval task such as ACTIVITYNET captions [16] and MSR-VTT [38]. The V3C dataset [27] does not use captions or subtitles but video metadata. While video descriptions or summaries are high-level representations of a video, they do not reflect human style search queries befitting of a realistic retrieval task. To more easily compare and contrast the existing datasets, we have compiled the details of said video datasets in Table 1.

**Video Retrieval Methods** Text-video retrieval is a core research area in video-language understanding [1, 4, 7, 31, 33, 34, 40, 41]. There have been many proposed solutions focusing on: combining pre-extracted features from frozen text and vision encoders [41]; transferring contrastive pre-training architectures, such as CLIP [24], to the video domain [10, 21, 22, 39]; leveraging video specific features such as sparsity [1, 20]; and aligning/fusing modalities [5, 6, 13, 15, 35, 37, 43]. Emphasis has been placed on aligning visual content with text, yet videos also include valuable information from text metadata, audio, and embedded text. Recently, fusion models [5, 6] and binding models [11, 43] have been developed to unify representations across modalities. VAST [6] notably integrates vision, audio, and subtitle information for text-video retrieval, enabling higher semantic video understanding in a unified representation. Despite these advances, many of these models are trained on datasets not designed for single-modality retrieval, let alone retrieval across multiple modalities. As a result, these systems face challenges when handling more complex retrieval tasks requiring inference across modalities and languages. This highlights the need for more comprehensive collections to facilitate the development of more robust and adaptable video retrieval systems capable of addressing the diversity of events found in real-world multimedia content.

### 3. Video Collection

Two significant challenges in current video retrieval research are: (1) existing collections often reward models for scene description alone, without considering the broader event context, and (2) these collections are typically limited in both size and scope. In this section, we outline the video collection process for MULTIVENT 2.0, demonstrating how it addresses these challenges. We first review the development of MultiVENT 1.0, which initiated efforts to tackle the first issue, and then describe the expanded collection process for MULTIVENT 2.0, a larger and more di-

verse dataset designed for real-world applications. The final MULTIVENT 2.0 dataset comprises over 217,000 videos, split evenly into train and test collections.

#### 3.1. MultiVENT 1.0 Development and Limitations

The creation of MULTIVENT 1.0 involved several key steps: developing topics based on language- and country-specific current events, collecting relevant videos for each topic, and aligning these events with corresponding news articles. For topic selection, Sanders et al. [29] utilized Google Trends statistics from countries with the largest populations of speakers for each target language to identify visually salient current events. After filtering to ensure adequate online video coverage, this process produced 2,396 videos spanning 255 events.

This principled approach resulted in an important dataset, filling a unique role as a targeted, event-centric video retrieval collection. However, while MULTIVENT 1.0 was a significant step forward, its small size, especially the limited distractor set, limits its applicability for large-scale multimodal retrieval research. As shown in Table 1, this issue is common across many video retrieval tasks, and stands in stark contrast to text retrieval benchmarks, which typically consist of much larger collections. For example, the HC4 dataset used in the 2022 NEUCLIR TREC shared task contains over 6 million text documents [17]. This disparity in collection size underscores the need for a larger and more diverse video retrieval corpus.

#### 3.2. Expanded MultiVENT 2.0 Video Collection Process

To create an expanded collection for a more realistic, large-scale video retrieval task, we augment the original MULTIVENT dataset with videos from INTERNVID, a corpus containing more than seven million YouTube videos and over 760,000 hours of content [36]. InternVid covers a superset of target languages and event categories than those found in MULTIVENT 1.0. While not all INTERNVID categories are event-based, the corpus still includes a significant amount of event-centric content, particularly in the political and disaster domains.

For MULTIVENT 2.0, we extract a large subset of videos from INTERNVID, filtering out those longer than five minutes. This process yields approximately 40,000 videos for each of the MULTIVENT 1.0 target languages: Arabic, Chinese, English, Korean, and Russian. To introduce greater variability at test time, we also include a smaller set of videos from the Spanish and “Unknown” language categories, ensuring that systems handle both new primary languages and a long tail of low-resource languages.

After merging the expanded collection with MULTIVENT 1.0, we compiled a final dataset of over 217,000 videos, consisting of 108,500 videos for train-

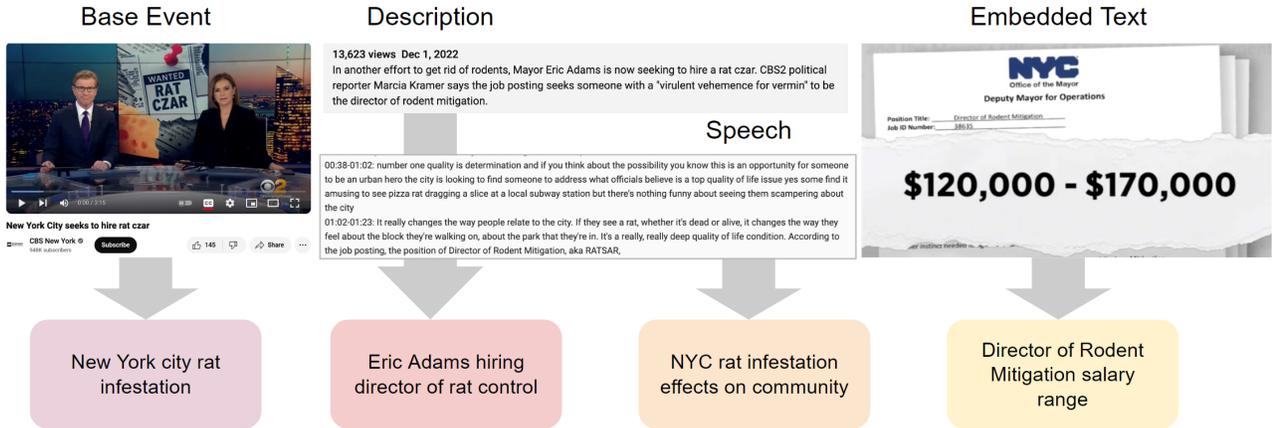


Figure 2. Query creation process for event-centric videos within our distractor collection from InternVid. Annotators first create a **Base Event** query based on the primary event depicted in the video. They then write up to three additional queries focusing on specific and unique aspects of the event: the **Description** query uses only information from the human-written text description, the **Speech** query relies on spoken content from the video, and the **Embedded Text** query utilizes text visible within the video frames.

ing (MULTIVENT TRAIN) and 109,800 for testing (MULTIVENT TEST). All videos from MULTIVENT 1.0 are solely found in the evaluation set, and any duplicates between the training and test collections were removed. Given the scale of the dataset, rapidly comparing models and system variations can be challenging. To mitigate this and aid model tuning, we offer a subset of 2,000 videos from the training set, referred to as MULTIVENT TRAIN-2K.

#### 4. Query Creation

As mentioned in Section 3, prior collections have focused more on matching descriptive aspects within videos. Accordingly, search queries have either been relatively short and vague (e.g., *a black and white horse runs around*) or derived directly from a video’s metadata, such as a YouTube description. MULTIVENT 1.0 largely followed this precedent, relying on the available metadata for each video. However, these approaches create a disconnect with modern text retrieval practices, where concise, search-engine-style English queries are increasingly used to match multilingual documents [17]. To bridge this gap, we develop a novel two-pronged approach for creating event-centric video queries: one that leverages additional fine-grained event annotations available with MULTIVENT 1.0 videos [29], and another that relies solely on the videos themselves.

For this task, we recruit professional linguists with expertise across the six primary languages targeted in MULTIVENT 2.0 to develop the queries. Each annotator underwent training via a tutorial task with detailed instructions. Afterward, we provide one-on-one feedback on their performance and closely monitor their initial annotations, offering additional guidance as needed. Finally, an annotation lead

conducts a quality control check to ensure the conciseness and specificity of each query/video pair.

##### 4.1. Updated Query Creation for MultiVENT 1.0

In MULTIVENT 1.0, the text descriptions accompanying each video were initially repurposed as queries. However, this dataset also provides a foundation for crafting more targeted, event-centric queries, as each of the 255 current events in the collection is aligned with an article in its original language, along with an English version if the source language is not English. Many of these English articles come from Wikipedia, where titles are typically concise and specific descriptions of the broader event. For instance, the article aligned with the event shown in Figure 1 is titled *2022 Lotus Garden China Telecom Building fire*. Annotators follow this style to manually create a base query for each current event in MULTIVENT 1.0, henceforth referred to as MULTIVENT BASE event queries.

While retrieving videos about overarching events is a valuable step forward, research in text information extraction (IE) has shown that individual events often encompass multiple distinct aspects, which are unlikely to be fully captured in a single piece of media. Thus, to further challenge retrieval models, annotators develop queries focused on specific, unique aspects of each event. This is facilitated by leveraging annotations from MULTIVENT-GROUNDED [30], a collection of 1,200 MULTIVENT 1.0 videos containing fine-grained event annotations. These annotations, guided by adapted FRAMENET event templates, include text description spans, video time intervals, and spatial bounding boxes at the frame level, and address IE-based questions such as *Where did the disaster occur?*, *Who was affected by the disaster?*, and *Who responded to the dis-*

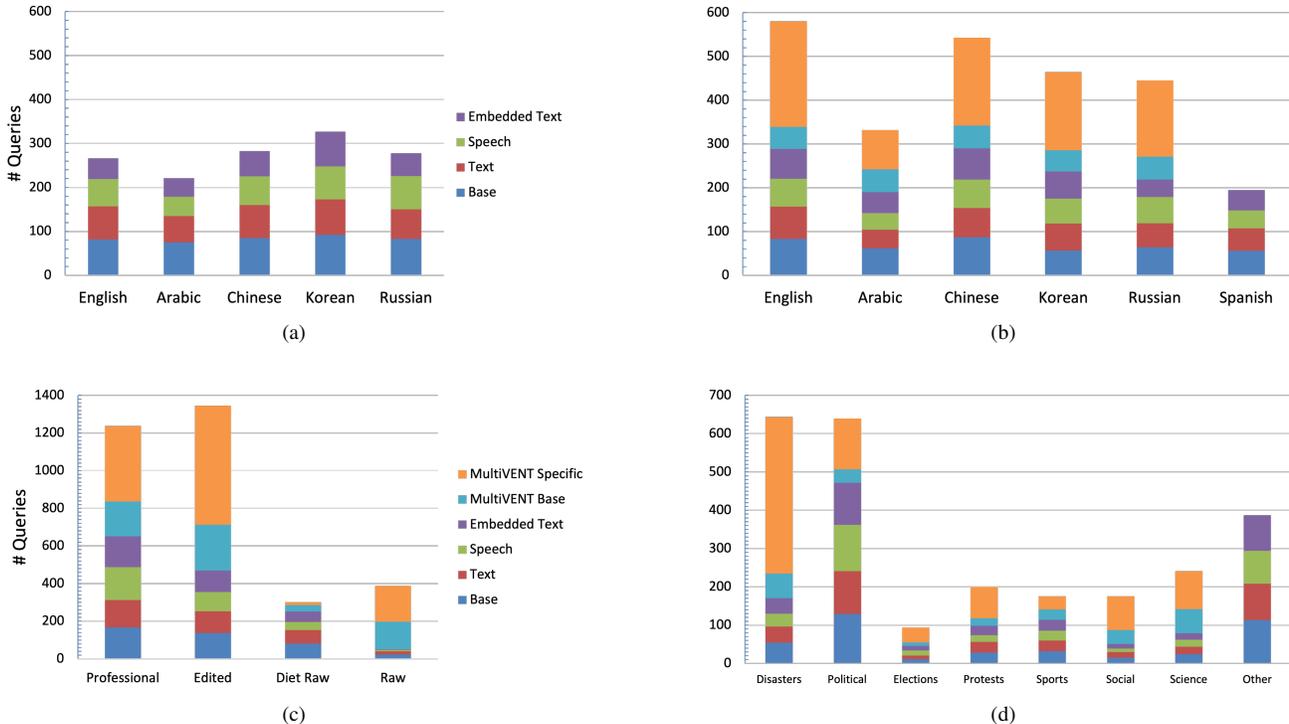


Figure 3. Breakdowns of the number of queries mapped to relevant videos. Figure 3a shows that MULTIVENT TRAIN contains queries targeting the five primary languages from MULTIVENT 1.0. On the other hand, in Figure 3b we see that MULTIVENT TEST adds queries targeting Spanish events to challenge systems’ multilingual robustness. Figure 3c shows that MULTIVENT 2.0 targets videos ranging from professional news broadcasts to raw first-person footage of events. Finally, as seen in Figure 3d, events in MULTIVENT 2.0 generally map to the same categories as MultiVENT 1.0, with a long tail of infrequent event types.

aster? Based on these annotations, annotators create 884 MULTIVENT SPECIFIC queries, each targeting videos that highlight distinct aspects of the corresponding event.

## 4.2. InternVid Query Creation

The previous section focuses on the retrieval of MULTIVENT 1.0 videos, despite the majority of our collection coming from InternVid. This raises the risk that part of the task could devolve into a binary classification problem, where simply distinguishing between MultiVENT and InternVid videos significantly reduces the task’s complexity. To mitigate this, we tasked annotators with writing queries targeting new events found within the larger InternVid collection. Since there was no guarantee that InternVid videos contained event-based content, we focused annotation efforts on videos from the News & Politics and Sports categories. Annotators were first asked to confirm whether a video contained event-based content. If confirmed, they write a **Base Event** query in the same style as those described in Section 4.1. Additionally, annotators searched the web for an article related to the event.

Next, annotators develop up to three additional queries, each focusing on specific and unique aspects of the same

event. For these queries, annotators were asked to rely on partial information about a video from a single modality: the **Description** query uses only the human-written YouTube description of the video; the **Speech** query uses any spoken content from the video, with Whisper-generated speech transcripts provided for assistance [25]; and the **Embedded Text** query focused solely on text directly visible in the video frames, with output provided from a multilingual OCR system [9]. Note that not all videos contain useful or unique information across all modalities; in such cases, no query is created. This process was applied to both the train and test collections, resulting in 1,417 additional test queries and 1,375 train queries.

## 4.3. Query and Video Breakdowns

To gain insight into the challenges posed by different types of video/query pairs, annotators are asked several additional questions about the videos during the query creation process. First, we ask them to confirm the video’s primary language. While INTERNVID ostensibly provides this information, the labels are sometimes incorrect; this is particularly common in English-labeled videos, where approximately 20% of the videos were in other languages. Fig-

Figure 3a shows the distribution of queries targeting videos from each of the five original MULTIVENT 1.0 languages within MULTIVENT TRAIN, while Figure 3b illustrates the same for MULTIVENT TEST. Notably, we add Spanish as an additional primary language in the test collection to evaluate the multilingual robustness of video retrieval systems.

Additionally, linguists group events into one of the seven event types defined in MULTIVENT 1.0: natural disasters (**Disasters**), political elections (**Elections**), **Protests**, other political developments (**Political**), sporting events (**Sports**), social events (**Social Events**), and scientific or technological discoveries (**Science**) [29]. An **Other** category is also provided for events that did not fit into these types. Figure 3d depicts the event type breakdown. Two notable aspects include the relative abundance of MultiVENT-specific queries for disaster events, largely due to MULTIVENT-GROUNDED annotations prioritizing this category [30], and the predominance of Political events, which reflects the news and politics videos comprising the majority of event-based content within the InternVid collection. The Other category includes diseases, police incidents, and man-made disasters, among other infrequent event types.

Finally, we ask annotators to categorize videos into three general types: **Professional** news broadcasts, i.e., videos with reporters and/or traditional news chyron; non-professional **Edited** videos, i.e., videos featuring multiple spliced clips, visual effects, or superimposed graphics; and **Raw** footage, i.e., single-stream videos of events as they happen, typically captured on a mobile device. Given the widespread availability of video editing software, it is very easy for more people to make light edits prior to uploading. To differentiate between this and "true" raw content, we add a fourth category, **Diet Raw**, for single-stream videos with minimal text and speech overlays. Figure 3c shows the distribution of queries for relevant videos across these types. From this, we can see that true raw content is relatively scarce in InternVid, reinforcing the importance of the original MultiVENT collection. Note that many queries targeting MultiVENT 1.0 events are mapped to multiple relevant videos across different video types.

## 5. Relevance Judgment Annotation

After the initial query creation process, we have a total of 6,068 annotated video/query pairs for MULTIVENT-TEST. In many cases, queries are only mapped to a single relevant video, and all other non-judged videos are assumed to be not relevant. This presents a challenge because, despite the targeted nature of our queries, it is likely that the distractor set contains additional relevant videos, given the size and diversity of the test collection. Furthermore, some videos may be partially relevant to a query, or relevant to some parts but not all, while others' relevance may be unclear due to a lack of context or inherent video ambiguity.

Rank	Very Relevant	Somewhat Relevant	Not Relevant
1	43%	13%	44%
2	28%	15%	57%
3	28%	19%	53%
4	20%	21%	59%
5	18%	18%	65%
6	14%	18%	69%
7	14%	17%	70%
8	10%	11%	78%
9	10%	15%	75%
10	9%	13%	79%

Table 2. Percentage of initially un-judged video/query pairs re-annotated as not relevant, somewhat relevant, and very relevant. The *Rank* indicates the ranking of a video for the corresponding query, as judged by our best baseline model.

To address this, we re-train the same pool of professional linguists from the query creation task in Section 4 to now judge the relevance of previously unseen query/video pairs. For each candidate video, annotators are asked to classify it as *not relevant*, *possibly relevant*, *partially relevant*, or *very relevant* to a query. Given the potentially vast scope of this annotation task, we limit the videos to those ranked in the top 10 by multilingual CLIP (MCLIP), a pre-trained vision-language model [2] and the strongest-performing model on MULTIVENT 1.0 [29]. After removing video/query pairs already judged, we prioritize judging the highest-ranked videos, and streamline the process by grouping query/video pairs where the same video was ranked in the top 10 for multiple queries. Through this process, we collect an additional 4,396 gold relevance judgments. For evaluation purposes, we condense the middle two categories into a single category labeled *somewhat relevant*.

Given these graded relevance categories, it follows that any video judged as very relevant for a **MultiVENT Base** query must be at least somewhat relevant for all associated **MultiVENT Specific** queries. To ensure systems receive appropriate credit, we apply this logic to add an additional 5,653 silver relevance judgments. This results in a final updated set of 16,116 judged video/query pairs. Prior to this annotation effort, on average only 16% of the videos ranked in the top 10 by multilingual CLIP had an associated relevance judgment (*Judged@10*); with the added judgments, the *Judged@10* has now increased to 39%. For evaluation metrics that account for graded relevance, such as normalized Discounted Cumulative Gain (nDCG), we apply a 0-1-3 scale: *Very relevant* video/query pairs are scored as 3, *Somewhat Relevant* pairs as 1, and *Not Relevant* or non-judged pairs as 0.

## 6. Baselines and Results

Jointly pre-trained vision-language models (VLMs) have recently achieved state-of-the-art results on various down-

Modality	Model	MSR-VTT	MultiVENT 2.0				
		R@10	R@10	R@100	MRR	mAP	nDCG@10
Vision	mCLIP	0.827	0.333	0.603	0.429	0.261	0.303
OCR	ICDAR → mCLIP	-	0.227	0.374	0.363	0.166	0.217
Speech	Whisper → mCLIP	-	0.290	0.450	0.417	0.212	0.267
Text	Description → mCLIP	-	0.293	0.491	0.445	0.228	0.284
Vision	InternVideo2.0	<b>0.851</b>	0.004	0.018	0.018	0.003	0.005
All	VAST	0.739	0.118	0.118	0.198	0.080	0.116
All	LanguageBind	0.787	<b>0.355</b>	<b>0.620</b>	<b>0.443</b>	<b>0.283</b>	<b>0.324</b>

Table 3. Performance of pre-trained multimodal benchmark models and single-modality pipeline systems. We can see that prior state-of-the-art video retrieval systems struggle significantly with the increase in query complexity and collection size/diversity of MultiVENT 2.0. VAST is run with vision input only on MultiVENT 2.0, as that achieves highest performance on this task.

stream video understanding tasks. These models typically consist of separate single-modality encoders, with a mechanism for fusing multimodal embeddings during training. In this section, we consider several prominent VLMs as baselines for video retrieval. Each of the following VLMs utilize a version of CLIP for the vision component, with the notable exception of InternVid using a Vision Transformer (ViT) [8].

- **VALOR** combines three encoders for single-modality representations with a decoder designed for multimodal text generation [5].
- **VAST** utilizes omni-modality pretraining to simultaneously fine-tune text, audio, and image encoders to enhance cross-modality learning [6].
- **InternVid 2** employs two expert models for video token-level unmasking during training and currently is state-of-the-art for zero-shot retrieval on MSR-VTT [37].<sup>3</sup>
- **Language-Bind** leverages a CLIP-based encoder for all non-text modalities and has achieved state-of-the-art results on several video-understanding tasks [43].

Beyond VLMs, we also evaluate several baseline pipeline approaches that utilize information from a single modality. In all cases, these system leverages mCLIP’s multimodal embedding space to efficiently compare model outputs to queries.

- **mCLIP**: We first extract ten keyframes by detecting significant scene changes within a video and then extracting the midpoint frame.<sup>4</sup> The frames are passed through mCLIP’s image encoder, and the pooled embeddings are then compared to analogous query embeddings.
- **ICDAR OCR**: Using the same frames from the vision baseline, we extract all visible embedded text using a state-of-the-art multilingual OCR system [9], and pass the

text output through mCLIP’s text encoder.

- **Whisper OCR**: We utilize WHISPER, a robust automatic speech recognition system [25], to transcribe each video’s audio track, and the resulting text is passed through mCLIP’s text encoder.
- **Description**: This baseline embeds each video’s human-written description through mCLIP’s text encoder.

Table 3 presents the performance of these baselines on MultiVENT 2.0, as well as results on MSR-VTT, a standard benchmark video retrieval datasets. A key observation is that while VLMs excel on prior collections, most perform poorly on our task. This is likely due to two main factors: the length of our videos, as VLMs are typically trained on short video segments; and the significant domain mismatch between simple visual concepts and complex event-based natural language, the latter being much harder to map directly to visual features. This second issue is compounded by some MultiVENT 2.0 queries targeting non-visual aspects of the video.

Another notable takeaway is that single-modality pipeline systems show promise, though there is considerable room for improvement. Table 4 highlights how different aspects of the problem can be better addressed by specific modalities. The distinction is particularly evident in the breakdown of video types in Table 4b: speech is highly effective for retrieving professional news broadcasts, yet it, along with embedded text, prove to be almost useless with raw content. In contrast, human-written descriptions and visual content are more robust across video types, with descriptions being the most useful for retrieving raw footage.

Based on these findings, a key question moving forward is whether to allow the use of text metadata. This may be appropriate in certain settings, but given the rapid rise of online visual content, it is unrealistic to assume that clean, human-written descriptions of a video will always be available. This is especially true for raw, user-generated con-

<sup>3</sup><https://paperswithcode.com/sota/zero-shot-video-retrieval-on-msr-vtt>

<sup>4</sup><https://www.senedetect.com>

Modality	Model	Performance by Language					
		Arabic	Chinese	English	Korean	Russian	Spanish
Vision	mCLIP	<b>0.323</b>	0.074	<b>0.331</b>	0.070	<b>0.247</b>	0.272
OCR	ICDAR → mCLIP	0.320	0.083	0.171	0.107	0.104	0.289
Speech	Whisper → mCLIP	0.232	<b>0.130</b>	0.203	0.114	0.224	<b>0.320</b>
Text	Description → mCLIP	0.246	0.127	0.199	<b>0.137</b>	0.242	0.281

(a)

Modality	Model	Performance by Video Type			
		Professional	Edited	Diet Raw	True Raw
Vision	mCLIP	0.201	<b>0.276</b>	<b>0.152</b>	0.198
OCR	ICDAR → mCLIP	0.184	0.176	0.094	0.006
Speech	Whisper → mCLIP	<b>0.318</b>	0.154	0.011	0.021
Text	Description → mCLIP	0.206	0.245	0.142	<b>0.307</b>

(b)

Modality	Model	Performance by Query Type					
		MultiVENT-Base	Multivent-Specific	Base	Text	Speech	Embedded Text
Vision	mCLIP	<b>0.411</b>	0.331	<b>0.294</b>	<b>0.294</b>	0.227	0.232
OCR	ICDAR → mCLIP	0.254	0.198	0.219	0.232	0.162	<b>0.245</b>
Speech	Whisper → mCLIP	0.268	0.225	0.282	0.256	<b>0.306</b>	0.230
Text	Description → mCLIP	0.403	<b>0.359</b>	0.219	0.282	0.160	0.153

(c)

Modality	Model	Performance by Event Type							
		Disaster	Election	Protest	Political	Sports	Social	Science	Other
Vision	mCLIP	0.334	<b>0.215</b>	<b>0.310</b>	<b>0.334</b>	<b>0.441</b>	<b>0.486</b>	<b>0.373</b>	0.259
OCR	ICDAR → mCLIP	0.208	0.153	0.244	0.186	0.284	0.287	0.228	0.218
Speech	Whisper → mCLIP	0.242	0.111	0.221	0.290	0.267	0.214	0.320	<b>0.321</b>
Text	Description → mCLIP	<b>0.343</b>	0.173	0.260	0.227	0.335	0.443	0.302	0.211

(d)

Table 4. **Breakdown of single-modality pipeline results.** Table 4a shows that Chinese and Korean videos are particularly challenging, especially for vision-only models. From Table 4b, extracted text is most beneficial for professional news broadcasts, while real-time raw footage remains challenging for all direct video content. Table 4c demonstrates that queries based on specific modalities are generally best handled by the corresponding modality. Finally, from Table 4d, visual content is most effective for most MultiVENT 1.0 event types.

tent, which is often posted in real-time as events unfold. As a result, we propose two versions of the final evaluation task for future systems: one where descriptions are disallowed (MULTIVENT TEST-NODESC) and one where they are permitted (MULTIVENT TEST-DESC).

## 7. Conclusion

In this paper, we introduce MULTIVENT 2.0, a massive multilingual, event-centric video retrieval collection that significantly broadens the scope and complexity of multimodal retrieval tasks. With over 218,000 videos and 3,900 queries targeting a diverse range of world events across six languages, MULTIVENT 2.0 poses substantial new challenges for vision-language models. Our results indicate that while pre-trained models such as Intern-

Vid and VAST have achieved impressive performance on prior retrieval collections, they face considerable difficulties in this setting, largely due to the domain mismatch between interpreting simpler visual concepts and the complex multimodal nature of event-based content. Pipeline systems incorporating specialized single-modality models show promise for subsets of the task, but these remain insufficient in isolation. In addition, retrieving raw video footage without the accompanied human-written text descriptions is extremely challenging, a limitation especially relevant in real-time event scenarios. Our results highlight the need for more robust multimodal retrieval systems, as exciting visual understanding and generation tasks cannot be leveraged in practice without first being able to accurately and efficiently triage relevant visual content.

## References

- [1] Meng Cao, Haoran Tang, Jinfa Huang, Peng Jin, Can Zhang, Ruyang Liu, Long Chen, Xiaodan Liang, Li Yuan, and Ge Li. Rap: Efficient text-video retrieval with sparse-and-correlated adapter, 2024. 3
- [2] Fredrik Carlsson, Philipp Eisen, Faton Rekathati, and Magnus Sahlgren. Cross-lingual and multilingual clip. In *Proceedings of the Language Resources and Evaluation Conference*, pages 6848–6854, Marseille, France, 2022. European Language Resources Association. 6
- [3] David Chen and William 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, Portland, Oregon, USA, 2011. Association for Computational Linguistics. 2
- [4] Shizhe Chen, Yida Zhao, Qin Jin, and Qi Wu. Fine-grained video-text retrieval with hierarchical graph reasoning, 2020. 3
- [5] Sihan Chen, Xingjian He, Longteng Guo, Xinxin Zhu, Weinong Wang, Jinhui Tang, and Jing Liu. Valor: Vision-audio-language omni-perception pretraining model and dataset, 2023. 2, 3, 7
- [6] Sihan Chen, Handong Li, Qunbo Wang, Zijia Zhao, Mingzhen Sun, Xinxin Zhu, and Jing Liu. Vast: A vision-audio-subtitle-text omni-modality foundation model and dataset, 2023. 3, 7
- [7] Ioana Croitoru, Simion-Vlad Bogolin, Marius Lordeanu, Hailin Jin, Andrew Zisserman, Samuel Albanie, and Yang Liu. Teachtext: Crossmodal generalized distillation for text-video retrieval, 2021. 3
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. 7
- [9] David Etter, Cameron Carpenter, and Nolan King. A hybrid model for multilingual ocr. In *Document Analysis and Recognition - ICDAR 2023*, pages 467–483, Cham, 2023. Springer Nature Switzerland. 5, 7
- [10] Han Fang, Pengfei Xiong, Luhui Xu, and Yu Chen. Clip2video: Mastering video-text retrieval via image clip, 2021. 3
- [11] Rohit Girdhar, Alaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *CVPR*, 2023. 3
- [12] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language, 2017. 3
- [13] Jingjia Huang, Yinan Li, Jiashi Feng, Xinglong Wu, Xiaoshuai Sun, and Rongrong Ji. Clover: Towards a unified video-language alignment and fusion model, 2022. 3
- [14] Po-Yao Huang, Mandela Patrick, Junjie Hu, Graham Neubig, Florian Metze, and Alexander Hauptmann. Multilingual multimodal pre-training for zero-shot cross-lingual transfer of vision-language models, 2021. 2
- [15] Peng Jin, Jinfa Huang, Fenglin Liu, Xian Wu, Shen Ge, Guoli Song, David A. Clifton, and Jie Chen. Expectation-maximization contrastive learning for compact video-and-language representations. In *Advances in Neural Information Processing Systems*, pages 30291–30306, 2022. 3
- [16] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos, 2017. 3
- [17] Dawn Lawrie, Sean MacAvaney, James Mayfield, Paul McNamee, Douglas W. Oard, Luca Soldaini, and Eugene Yang. Overview of the trec 2022 neuclir track, 2023. 1, 2, 3, 4
- [18] Jie Lei, Licheng Yu, Tamara L. Berg, and Mohit Bansal. Tvr: A large-scale dataset for video-subtitle moment retrieval, 2020. 2
- [19] Jie Lei, Tamara L. Berg, and Mohit Bansal. Mtvr: Multilingual moment retrieval in videos, 2021. 2
- [20] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling, 2021. 3
- [21] Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An empirical study of clip for end to end video clip retrieval, 2021. 3
- [22] Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. X-clip: End-to-end multi-grained contrastive learning for video-text retrieval. In *Proceedings of the 30th ACM International Conference on Multimedia*, page 638–647, New York, NY, USA, 2022. Association for Computing Machinery. 3
- [23] 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, 2019. 2
- [24] 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, 2021. 3
- [25] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022. 5, 7
- [26] Anna Rohrbach, Atousa Torabi, Marcus Rohrbach, Niket Tandon, Christopher Pal, Hugo Larochelle, Aaron Courville, and Bernt Schiele. Movie description, 2016. 2
- [27] Luca Rossetto, Heiko Schuldt, George Awad, and Asad A. Butt. V3c - a research video collection, 2018. 3
- [28] Kate Sanders and Benjamin Van Durme. A survey of video datasets for grounded event understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7314–7327, 2024. 2
- [29] Kate Sanders, David Etter, Reno Kriz, and Benjamin Van Durme. Multivent: Multilingual videos of events with aligned natural text, 2023. 1, 3, 4, 6
- [30] Kate Sanders, Reno Kriz, David Etter, Hannah Recknor, Alexander Martin, Cameron Carpenter, Jingyang Lin, and Benjamin Van Durme. Grounding partially-defined events in multimodal data, 2024. 4, 6

- [31] Haoran Tang, Meng Cao, Jinfa Huang, Ruyang Liu, Peng Jin, Ge Li, and Xiaodan Liang. Muse: Mamba is efficient multi-scale learner for text-video retrieval, 2024. [3](#)
- [32] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex: A large-scale, high-quality multilingual dataset for video-and-language research, 2020. [2](#)
- [33] Xiaohan Wang, Linchao Zhu, and Yi Yang. T2vlad: Global-local sequence alignment for text-video retrieval, 2021. [3](#)
- [34] Yimu Wang and Peng Shi. Video-text retrieval by supervised sparse multi-grained learning, 2023. [3](#)
- [35] Yi Wang, Kunchang Li, Yizhuo Li, Yanan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, Sen Xing, Guo Chen, Juntong Pan, Jiashuo Yu, Yali Wang, Limin Wang, and Yu Qiao. Internvideo: General video foundation models via generative and discriminative learning, 2022. [3](#)
- [36] Yi Wang, Yanan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, Conghui He, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Yu Qiao. Internvid: A large-scale video-text dataset for multimodal understanding and generation, 2024. [3](#)
- [37] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yanan He, Chenting Wang, Guo Chen, Baoqi Pei, Ziang Yan, Rongkun Zheng, Jilan Xu, Zun Wang, Yansong Shi, Tianxiang Jiang, Songze Li, Hongjie Zhang, Yifei Huang, Yu Qiao, Yali Wang, and Limin Wang. Internvideo2: Scaling foundation models for multimodal video understanding, 2024. [3](#), [7](#)
- [38] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5288–5296, 2016. [1](#), [3](#)
- [39] Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. Clip-vip: Adapting pre-trained image-text model to video-language representation alignment, 2023. [3](#)
- [40] Jianwei Yang, Yonatan Bisk, and Jianfeng Gao. Taco: Token-aware cascade contrastive learning for video-text alignment, 2021. [3](#)
- [41] Youngjae Yu, Jongseok Kim, and Gunhee Kim. A joint sequence fusion model for video question answering and retrieval, 2018. [3](#)
- [42] 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*, 2018. [2](#)
- [43] Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiayi Cui, Hongfa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, Wancai Zhang, Zhifeng Li, Wei Liu, and Li Yuan. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment, 2024. [3](#), [7](#)

# MultiVENT 2.0: A Massive Multilingual Benchmark for Event-Centric Video Retrieval

## Supplementary Material

### 8. Query Creation - Annotation Guidelines

In this section, we provide the guidelines given to our professional linguists for the task of writing queries targeting aspects of new events found within our larger collection sampled from InternVid, described in Section 4. This annotation task was set up and run in Label Studio, a flexible data labeling platform that allows for multimodal annotations.<sup>5</sup> When writing queries for MULTIVENT-TRAIN and MULTIVENT-TEST, we ran two iterations of this task: the first looked at a random sample of 200 Arabic, Chinese, English, Korean, and Russian videos. For MULTIVENT-TEST we added Spanish as a sixth language. We focused on annotating InternVid videos from the News & Politics and Sports domains, as these were the categories with the most event-focused content.

After this initial round, there was not sufficient raw event-focused videos. To address this, we again sampled 200 videos per language from InternVid, this time filtering out any videos with standard resolutions, and continued until we had at least 20 annotated videos for each language/video type pair.

#### 8.1. Task Introduction

The goal of this task is to annotate videos based on event templates to identify relevant information/entities across multiple modalities (text descriptions, video footage, text in videos, and audio/speech). These annotations were used to create queries based on the MULTIVENT 1.0 dataset and the events it contains. For the retrieval task to be meaningfully challenging, the set of documents from which systems are retrieving needs to be quite large (in the hundreds of thousands). As such, researchers have supplemented the MULTIVENT dataset with videos from a larger collection of videos, INTERNVID, and have begun writing queries associated with videos from that set.

In this annotation task, language experts will write queries for INTERNVID videos, with a focus on queries that come from specific modalities. For each annotation, you will be presented with a video (from YouTube) and its description (if present). Your task will be to watch the video, read the description, and write up to 4 English search queries based on information in the video, video description, video audio, and video text.

<sup>5</sup><https://labelstud.io>

### 8.2. Video Relevance

**What is an event?** We broadly define an “event” as some set of properties and their changes (or lack of change) over time. A single, distinct event can only exist at one set location(s) across one set time interval. For this task, we generally constrain an event to something that occurs naturally, that is, not a staged set of predetermined actions (think movies, tutorial videos, etc.). We generally want to write queries that pertain to the main events depicted in the videos.

**When should I skip a video?** If the video falls into one of the below categories, select *Not Relevant* and submit the annotation.

- If it is hard to tell what the main event in the video is.
- If there is minimal event-relevant information present in the video (for example, a black screen, or a slideshow of irrelevant images).
- If the video is staged or otherwise not natural (movies, commercials, tutorial videos, etc.).
- If the video is not mostly in the language you are annotating for (this applies to the description as well as the audio content).
- If you are uncomfortable watching the video. Please do not annotate any videos that you don’t want to watch!

### 8.3. Video Classification and Video Type

Select an event type that best describes the event in your video. If none of the listed event types apply, select *Other event type*.

- Disaster (fire, earthquake, hurricane, etc.)
- Political Election
- Political Protest
- Other Political Development (sanctions, treaty, leader death, etc.)
- Social Event (festival, convention, celebration, etc.)
- Sporting Event
- Technical/Scientific Launch or Discovery
- Other Event Type

Then, indicate the type of video:

- Professional news broadcast
- Non-professional edited footage
- Mostly raw footage
- Raw footage

The guidelines for video type classification are:

- The video is **professional footage** if you can imagine the clip appearing on broadcast TV. For example, there is an anchor/reporter, there are news logos and text banners with headlines, etc.
- The video is **non-professional edited footage** if it is not professional AND there is superimposed OCR/graphics, other visual special effects have been applied, or there are multiple clips spliced together. For example, many YouTube streamers’ videos are “edited non-professional.”
- The video is **mostly raw footage** if the video is a mostly continuous raw stream of content but has minor additional content (i.e., one of speech/music overlay, text overlay, or minimal scene cuts).
- The video is **true raw footage** if it is a continuous, unedited stream of footage.

## 8.4. Query Writing

Tips for good query writing:

- All of the queries should be in English.
- Write queries in the style you would use for a search engine, e.g., Google. The queries should be short and do not need to be written in complete sentences.
- Only write queries based on information that would help you learn more about the event in question. There may be video, audio, or text content that is not relevant to the main event topic—these irrelevant details should not be the basis for your queries.
- The queries should be unique from one another. If you cannot write a query for unique information in a particular modality, check the relevant box for why you left the query blank (see below).
- You can use supplemental internet searches to clarify information contained in the video description, e.g., to check spelling for a person’s name.
- There are multiple “right answers” for each video’s queries.

Each query has an accompanying text box; write your query and hit Enter to submit that query. Once submitted, the query will appear in a green text box.

### 8.4.1 Base Event Query

Watch the video. You may also read the title/description to get a better sense of what the video is depicting. Then, write one query based on the **primary topic/event** depicted in the video. This should include general information about what happened. The goal is to capture specific events. Each query should be approximately 3-7 words. An example of this question is provided in Figure 4.

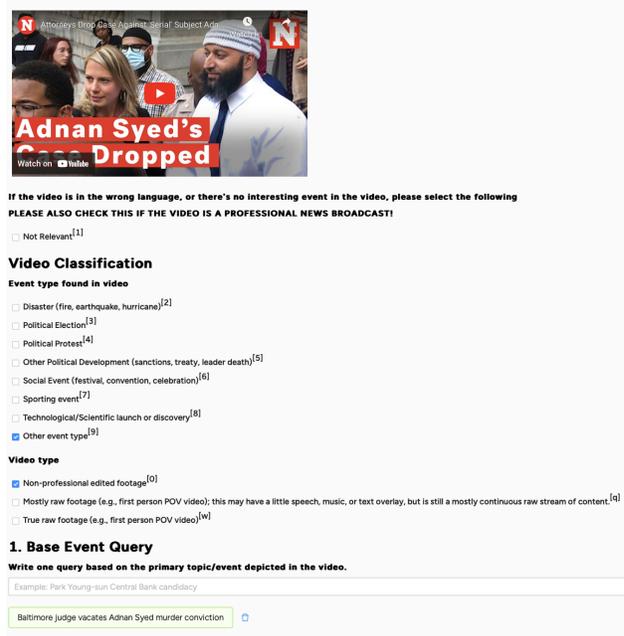


Figure 4. Example and Label Studio interface for video classification annotations and base query creation i.e., the phrase that represents the “Wikipedia title” of a current event.

### 8.4.2 Text-based Query

Write one query based on event-relevant details from the video description. This query should cover information that is not already addressed by previous queries. The query should be approximately 3-7 words. An example of this question and its associated interface is provided in Figure 5a.

If no unique query can be written, select an option explaining why:

- No event-relevant information in the description.
- Cannot write a unique query based on the description.

### 8.4.3 Audio-based Query

Write one query based on event-relevant details from the video audio/speech. This query should cover information that is not already addressed by previous queries. The query should be approximately 3-7 words.

To aid annotation, transcribed speech has been provided (in the original language or translated, depending on the project you are working on). Note that these transcriptions, especially once translated, can have errors. Listen to the audio if you are a language expert and try to confirm spellings of names, places, etc. before finalizing your queries. An example of this question and its associated interface is provided in Figure 5b.

## 2. Text-based Query

### Title

Attorneys Drop Case Against 'Serial' Subject Adnan Syed: 'Justice Is Done'

### Description

Baltimore prosecutors dismissed charges against Adnan Syed Tuesday, citing new evidence. Syed was convicted in 2000 for the killing of ex-girlfriend Hae Min Lee. His case was later popularized in hit podcast 'Serial,' which raised questions on his conviction. Subscribe on YouTube: <https://bit.ly/2Jpnlb>

-----  
Facebook: <https://www.facebook.com/Newsweek>  
Twitter: <https://twitter.com/Newsweek>  
Instagram: <https://www.instagram.com/Newsweek>  
<https://www.newsweek.com>

**Write a new and unique query based on event-relevant details from the video title and description.**

Example: Park Young-sun press conference

Adnan Syed 2000 murder conviction

**2a. If no query can be written, select one of these options for why:**

- No event-relevant information in the description.<sup>[a]</sup>
- Cannot write a unique query based on the description.<sup>[t]</sup>

(a) Query creation using only the associated text description

## 3. Audio-based Query

### Transcribed speech

00:00-00:25: This morning, I instructed my office to dismiss the criminal case against Adnan Saeed following the completion of a second round of touch DNA testing of items that were never tested before. Those items include skirt, pantyhose, shoes, and jacket of Miss Hae Min Lee.  
00:26-00:47: Although no DNA was recovered from the skirt, the pantyhose, or jacket swabs, there was a DNA mixture of multiple contributors on both Ms. Lee's shoes. The same multiple contributor for both of Ms. Lee's shoes. And most compellingly, Adnan Saeed  
00:48-01:14: His DNA was excluded. We stand ready and willing to provide whatever counseling or support services that may be needed for that family who has had to relive an unimaginable nightmare over and over again. Equally heartbreaking is the pain and the sacrifice and the trauma  
01:15-01:41: that has been imposed not just on that family, but Adnan and his family, who together spent 23 years in prison for a crime as a result of a wrongful conviction. As the administrator of the criminal justice system, it's my duty to ensure that justice is not delayed, justice is never denied, but justice be done.  
01:41-01:59: Today, justice is done. And that means today, tomorrow, and until my administration ends, we will continue to utilize every available resource to prosecute whoever is responsible for the death of Hae Min Lee.

**Write a new and unique query based on event-relevant details from the audio.**

Example: XX

New DNA testing in decades-old Baltimore murder case

**3a. If no query can be written, select one of these options to explain why. A provided transcript was generated automatically and isn't perfect, so please listen to the full audio before deciding.**

- No event-relevant information in the audio.<sup>[a]</sup>
- Cannot write a unique query based on the audio.<sup>[a]</sup>

(b) Query creation process using only audio, with associated transcribed text for assistance

## 4. OCR-based Query

### Text found within the video (OCR)

00:03: office to dismiss the criminal || This morning, I instructed my || Baltimore City State's Attorney || Marilyn J. Mosby || Oct. 11,2022 || Baltimore, Maryland  
00:11: testing of items that were never || a second round of touch DNA || Associated Press/WMAR  
00:19: include a skirt, pantyhose, shoes || tested before. Those items  
00:27: Although no DNA was recovered || and jacket of Miss Hae Min Lee.  
00:35: on both Miss Lee's shoes. || mixture of multiple contributors  
00:43: for both of Miss Lee's shoes. || The same multiple contributor  
00:50: provide whatever counseling or || We stand ready and willing to  
00:58: needed for that family who has || support services that may be  
01:06: and the sacrifice and the trauma || Equally heartbreaking is the pain  
01:22: years in prison for a crime as a || family, who together spent 23  
01:30: criminal justice system, it's my || As the administrator of the  
01:38: delayed, justice is never denied, || duty to ensure that justice is not  
01:45: and until my administration || And that means today, tomorrow  
01:53: every available resource to || ends, we will continue to utilize  
02:01: Mickey Hutchings || Produced by || Newsweek || WMAR/AP || Visuals

**Write one new and unique query based on event-relevant details from OCR.**

Example: XX

New evidence in Hae Min Lee murder

**4a. If no query can be written, select one of these options to explain why. The provided OCR was generated automatically and isn't perfect, so please look at all text within the video before deciding.**

- No event-relevant information in the text within the video.<sup>[d]</sup>
- Cannot write a unique query based on the text within the video.<sup>[t]</sup>

(c) Query creation process using only embedded text, with associated OCR output for assistance

Figure 5. Example and Label Studio interface for writing queries targeting specific aspects of events. For each question, we ask annotators to only use the text description, audio, and embedded text, respectively.

If no unique query can be written, select an option explaining why:

- No event-relevant information in the audio.
- Cannot write a unique query based on the audio.

#### 8.4.4 OCR-based query

Write one query based on event-relevant details from text in the video footage (e.g., signs, captions, etc.). Read the text found within the video (*OCR*) section. This query should cover information that is not already addressed by previous queries. The query should be approximately 3-7 words.

To aid annotation, OCR transcripts have been provided (in the original language or translated, depending on the project you are working on). Note that these transcriptions, especially once translated, can have errors and omissions. Consult the video footage for text that may have been missed or misrendered by OCR and consider taking a screenshot and re-translating in Google Translate if you are not a language expert before finalizing your queries. An example of this question and its associated interface is provided in Figure 5c.

If no unique query can be written, select an option explaining why:

- No event-relevant information in the text within the video.
- Cannot write a unique query based on the text within the video.

## 9. Relevance Judgment - Annotation Guidelines

In this section, we provide the guidelines given to our professional linguists to judge the relevance of previously unseen query/video pairs, described in Section 5.

### 9.1. Introduction

The goal of the task is for models to differentiate, i.e., through a ranked list, between relevant and irrelevant documents with respect to each query. Using Google as an example use case, a user expects the content most relevant to their search to appear in the first page of results. Through our prior query creation, we have developed a set of over 3,900 unique queries. To make the retrieval task meaningfully difficult, we have also included a large set of distractor videos (videos for which we have not written queries) in the overall document set.

When evaluating models on the retrieval task, our metrics assume that any given unseen video will be *Not Relevant* for any given query. However, it is possible that there are videos in our distractor set that are relevant for some of our queries. We want to find as many relevant videos as possible for a query to reduce the possibility of accidentally labeling something not relevant that actually is (called a *false negative*). This means potentially increasing the number of

videos that are judged as relevant. This will help us more accurately judge systems.

Providing Relevance Judgment annotations involves assessing a query-video pair and determining if, for that query, the video is *Very Relevant*, *Somewhat Relevant*, or *Not Relevant*. The sections below provide more detailed guidance for this task.

### 9.2. Annotation Instructions

For each annotation task, you will be presented with a video (from YouTube), its title and description (if present), and extracted speech and embedded text. The text will be translated to English, if necessary (if you are a language expert and want to see the original title/description, open the video in a separate window in YouTube).

#### 9.2.1 Video Type

After watching the video, indicate the type of video—professional footage/news broadcast, non-professional edited footage, or raw footage (i.e., first person POV video). An example of this component is shown in Figure 6a.

The guidelines for video type classification are:

- The video is **professional footage** if you can imagine the clip appearing on broadcast TV. For example, there is an anchor/reporter, there are news logos and text banners with headlines, etc.
- The video is **non-professional edited footage** if it is not professional AND there is superimposed OCR/graphics, other visual special effects have been applied, or there are multiple clips spliced together. For example, many YouTube streamers' videos are "edited non-professional."
- The video is **mostly raw footage** if the video is a mostly continuous raw stream of content but has minor additional content (i.e., one of speech/music overlay, text overlay, or minimal scene cuts).
- The video is **true raw footage** if it is a continuous, unedited stream of footage.

Note: some news agencies/platforms release a mix of video types. A "non-professional edited" video may be from a professional news company, but not meet our criteria for professional footage.

#### 9.2.2 Relevance Judgments - Overview

Next, you will see a list of 1 or more English queries that you will judge that video against for relevance. For each query, you will select one of the following options:

- **Not Relevant**

- **Possibly Relevant**, i.e., the query mentions a clear event, and the video is possibly about that event
- **Partially Relevant**, i.e., the query mentions a clear event, the video is about that event, but does not show everything in the query
- **Very Relevant**

You will also have a text box to optionally leave comments (e.g., if you are not sure and/or escalating the judgment and want to briefly explain your reasoning).

### 9.2.3 About Our Queries

**Base event queries** When we wrote the base event queries, we tried to identify the main event being addressed by a given video. The model for this approach was MultiVENT base event queries, which were based on trending events that i.e. were covered by a Wikipedia article. This means that often the base event queries are broader than the videos they were written for, which may only cover specific aspects of the event in question. Some examples are below:

- A video covering loosening mask requirements in March 2022 in Canada might have the base event query *Covid-19 pandemic in Canada*.
- A video depicting a post-game interview with Chiefs quarterback Patrick Mahomes might have the base event query *Super Bowl LVIII*.
- A video of highlights from a presidential debate might have the base event query *2022 South Korean presidential election*.

For each of the query-video pairs above, the video would be Very Relevant for that query, even though it may not address the event as a whole. The base event queries may also name the event using details not explicitly mentioned in the video (for example, the year). However, if you are able to determine that the event described in the query is the same underlying event that the video is about, that video will be *Very Relevant*.

**Specific Queries** We wrote specific queries based on details found in videos. Specific queries will usually refer to the event in question as well as a specific aspect of that event. When you judge relevance for specific queries, you will want to consider (to the extent that they are decomposable) both the event described by the query and the specific aspect(s) or detail(s) about the event included in the query. Some examples are below:

- For a query *Black Sea drone incident damaged propeller*, you will want to consider 1) is the video about the Black Sea drone incident and, if yes, 2) does the video also include reference (in the image, description, speech, and/or embedded text) to a damaged propeller?

- For a query *Democratic Party response to Daejang-dong scandal*, you will want to consider 1) is the video about the Daejang-dong scandal, and 2) does the video also include content about the Democratic Party's response to the scandal?

If a specific query does not explicitly reference the *base event*, evaluate relevance based on each detail in the query. Some examples are below:

- For a query *Eric Adams hiring Director of Rat Control*, consider whether the video explicitly refers to Eric Adams (or New York Mayor) and to Director of Rat Control/Director of Rodent Mitigation.
  - If the video is about Eric Adams but unrelated to rat control, the video is Not Relevant.
  - If the video is about hiring the Director of Rat Control but does not refer to Eric Adams/the New York Mayor, the video is Partially Relevant.
- For a query *Canada eases entry restrictions September 2022*, consider whether the video is about updates to entry restrictions, in Canada, in September 2022.
  - Because the query points to a specific event, if the video is about Canadian entry restrictions not in September 2022, the video is Not Relevant.
  - If the video is about Canadian entry restrictions and the date/time period is not clear, the video is Partially Relevant.

### 9.2.4 Consulting Outside Sources

If you are unfamiliar with the event, people, places, etc. described by a given query, we encourage you to search for the event and/or entities online to learn more about the query's context. For named entities, this is particularly helpful, because the queries might refer to those entities in different ways than the videos do. For example, a query might say *Norilsk avalanche* and a video might refer to an *avalanche in the Krasnoyarsk region*; if you are not familiar with the region, an internet search may be required to determine that the video is relevant for that query.

When assessing the queries and videos for relevance, try to think conceptually about what a query is asking for, rather than looking for that exact wording in the video.

### 9.2.5 About Our Videos and Metadata

Most videos you will see for this task are from YouTube. You may also encounter some from Twitter. In addition to the video, you will also see various metadata for that video, shown in Figure 6b. For videos not originally in English, the metadata will have been translated into English.

- **Video Title + Video Description.** If you are a language



Please select one of the following options to escalate this task to a reviewer.

- Escalate<sup>[1]</sup>
- Video removed<sup>[2]</sup>

**Video type**

- Professional footage (news broadcast)<sup>[3]</sup>
- Non-professional edited footage<sup>[4]</sup>
- Mostly raw footage (e.g., first person POV video); this may have a little speech, music, or text overlay, but is still a mostly continuous raw stream of content.<sup>[5]</sup>
- True raw footage (e.g., first person POV video)<sup>[6]</sup>

(a) Example video and classification annotations

**This is the video title.**

Top secret spy planes reveal details about China's spy balloon

**This is the date the video was posted online.**

02-10-2023

**This is the video description.**

As tensions between the US and China escalate following the US downing and recovery of a suspected spy balloon, CNN's Will Ripley takes a look at the claims made by the respective countries and details surrounding the balloon and China's vast military surveillance program.  
#CNN #News

**Audio-based Information (Transcribed speech).**

Tonight, intelligence from top-secret U.S. spy planes is providing startling details about China's spy operations. We are learning that images taken of the spy balloons by high-altitude U-2 planes show in detail five antennas and huge solar panels, the kind of technology that gave the Chinese not just the ability to listen to U.S. communications and to pinpoint exactly where they were coming from, who they were coming from. All of this especially alarming as the balloon passed over sensitive military sites before the U.S. Air Force shot it down off the South Carolina coast. And of course, we now know there have been other balloons where that one came from. And these new details come as China is forcefully pushing back against President Biden's very public criticism of Chinese President Xi Jinping at the State of the Union. Out front now, Will Ripley, who tonight is in Taiwan. Fiery new rhetoric from China, escalating the suspected spy balloon scandal. Beijing blasting President Joe Biden for criticizing Chinese President Xi Jinping. Can you think of any other world leader in trade places with Xi Jinping? Not a joke. You think of any? Who would? I can't think of one. This man has enormous problems. China says Biden's remarks are highly irresponsible and violate basic diplomatic protocols. Problems complicated by a growing pile of evidence. Pieces of the downed balloon pulled from the sea off the Carolina coast, proof the Pentagon's China's weather balloon claim is nothing but hot air. A Chinese foreign ministry spokesman says, I have no knowledge about America's claim that this balloon is part of a fleet. I think it could be part of the information and public opinion war that the US is waging against China. The international community can see clearly who's the world's largest espionage and surveillance country. I can assure you this was not for civilian purposes. That is, we are 100%, clear about that. There's that. The US linked the balloon to a vast Chinese military surveillance program, a growing list of global balloon sightings and questions. The US believes many balloons are launched from China's Hainan Island, where a US spy plane made an emergency landing in 2001. China took three months to investigate before returning the plane in pieces. Now China is attacking the US for shooting down its balloon and sending the pieces to an FBI lab. The spokesman says the US insists on using force to attack Chinese unmanned civilian airships, which seriously violates international practice and sets a horrible precedent. US Defense Secretary Lloyd Austin asked for a phone call with China's defense minister. Beijing bluntly declined. And, Will, it's a precarious moment, and I know that you have been learning that Xi Jinping was directly involved. Yes. I'm going to ask you about the balloon, the balloon that was flying over the US. And I know that the balloon is not new. And I know that the balloon is not new in this program. And, in fact, balloons flying over the US spying from China is not new. But just how much did Xi Jinping know? Yeah, I mean, you're talking about more than 40 countries across five continents. According to a senior State Department official, these balloons have gathered intelligence. This program has been going on with President Xi's direct knowledge. However, lawmakers briefed on this have been told that this particular balloon at this particular time was a fake. And that the balloon, which was on display at the time, may have been dispatched without the knowledge of Xi and also without the knowledge of senior People's Liberation Army and Communist Party leaders, Aaron.

**OCR-based Information (Text found within the video).**

00:06: UKRAINIAN PRES. VOLODYMYR ZELENSKI ERIN BURNETT || CN.com || 7:56 PM ET || U.S.: CHINESE SPY BALLOON CAPABLE OF MONITORING U.S. COMMUNICATIONS || CNN || LIVE || NEW DEVELOPMENTS || Myrtle Beach, SC || Jonathan Snyder /LSM Saturday  
00:20: "THIS IS OUR EUROPE, ERIN BURNETT || " AS HE PUSHES FOR MORE MILITARY SUPPORT || NAS T -120.94 || HAD ANTENNAS, CAPABLE OF MONITORING U.S. COMMUNICATIONS || U.S.: IMAGES FROM U.S. SPY PLANES REVEALS CHINESE SPY BALLOON || Jonathan Snyder /LSM | Saturday  
00:33: OUTFRONT || FOR UKRAINE, IT'S A WAY HOME," ZELENSKY SAYS, REFERENCING UKRAINE  
00:47: BECAME AN EU CANDIDATE STATE LAST YEAR, BUT BECOMING AN OFFICIAL || NAS T -120.94 || OL  
01:00: TRIES THAT HAVE ALREADY PROVIDED WEAPONS AND MILITARY, BUT STRES ERIN BURNETT || 4:57 PM PT  
01:14: ICE MISSILES, AND MODERN FIGHTER, JETS TO PROTECT ITS SECURITY, WHICH ERIN BURNETT || AND VIOLATE BASIC DIPLOMATIC PROTOCOLS" || रश्मि  
01:27: "S TRIP TO BRUSSELS COMES A DAY AFTER HE MADE SURPRISE VISIT TO LO || US Navy  
01:41: MEANTIME, RUSSIA IS BOOSTING || DMATIC TOUR OF EUROPEAN CAPITALS || T -36.36 || S&P || U.S.: IMAGES FROM U.S. SPY PLANES REVEALS CHINESE SPY BALLOON  
01:54: ZELENSKY'S PLEAS, ACCORDING TO PUTIN ALLY AND FORMER RUSSIAN PRE ERIN BURNETT || DOW T -249.13 || | | Defense  
02:08: "YESTERDAY, OUR ENEMY BEGGE ERIN BURNETT || OF THE RUSSIAN SECURITY COUNCIL || T -36.36 || Taiwan || Costa Rica || Missouri

(b) Provided text extracted from the video

**Relevance Judgments**

Taking all available information into account, judge the video's relevance in relation to the following query. If there is no query, you can ignore the question or fill out not relevant

**Query 1**

Travel delays to China afte 2023 spy balloon controversy

- Not Relevant<sup>[7]</sup>
- Possibly Relevant, i.e., the query mentions a clear event, and the video is possibly about that event<sup>[8]</sup>
- Partially Relevant, i.e., the query mentions a clear event, the video is about that event, but does not show everything in the query<sup>[9]</sup>
- Very Relevant<sup>[0]</sup>

**Query 2**

Biden retaliation 2023 spy balloon incident

- Not Relevant<sup>[4]</sup>
- Possibly Relevant, i.e., the query mentions a clear event, and the video is possibly about that event<sup>[4]</sup>
- Partially Relevant, i.e., the query mentions a clear event, the video is about that event, but does not show everything in the query<sup>[4]</sup>
- Very Relevant<sup>[1]</sup>

(c) Queries for which the video is judged for relevance.

Figure 6. Example of the annotation process for judging relevance for previously unseen query/video pairs.

expert and want to see the title and description in the original language, open the video in another window.

- **Posted Date.** Note that while it is helpful to use this date as a clue for the event depicted by the video, the date the video was posted online will not necessarily correspond closely to the date of the event. Consider the date along with all other information available to you.
- **ASR (automated speech recognition) output.** Speech from the video has been automatically transcribed, then translated into English if necessary.
- **OCR (optical character recognition) output.** A sample of text from the video (i.e., captions, banners, signs) has been automatically transcribed, then translated into English if necessary. Not all video frames were considered for this OCR output—this means that there may be text that was not caught at all for transcription and translation. You can supplement this information by pausing the video on a frame that has text on it, taking a screenshot of the video frame, and pasting the image into Google Translate for images.

Each automated process described above has the potential for errors. This is especially true for the recognition and translation of proper nouns, e.g., names. Consider all of the information available in the video, metadata, and outside sources to judge whether details in the video correspond to details from the query.

## 9.2.6 Relevance Judgments - Detailed Guidance

A video is Relevant with respect to query if some content in the video addresses the information-seeking need defined by the query. How primary that relevant content is within the video (e.g., if it is the main topic of the videos vs. if it is mentioned in passing) does not affect the video's relevance for the query.

Each query defines what kind of information we are looking for in the video (see discussion of Base Event queries and Specific queries above). Some videos will be Very Relevant (i.e., they completely address all aspects of a query), others will be Not Relevant (they do not address what the query is looking for). There are also two intermediate relevance judgments, Partially Relevant and Possibly Relevant. Example queries are provided in Figure 6c, and more details about each category are below:

**Very Relevant [3]** A video is *Very Relevant* for a query if it addresses all components of a query.

- For a base event query, a *Very Relevant* video is one about that event.
- For a specific query, a *Very Relevant* video is one that is about the event described by the query (if applicable) and includes the specific information/details in the query.

**Partially Relevant [2]** A video is *Partially Relevant* for a query if it addresses the general topic of the query but does not include the specific detail(s) the query is asking for.

- For a specific event query, a *Partially Relevant* video is one that covers the main event/topic in the query but is missing a specific detail included in the query.
- If you can confirm that a query is about an event distinct from the one in the query, the video is Not Relevant.

**Possibly Relevant [1]** A video is *Possibly Relevant* if it is unclear that the video is relevant for the query.

- For a base or specific event query, a *Possibly Relevant* video is one that may cover an aspect of the event, but it is not entirely clear that it does so (i.e., it is possible that the video is about the event/topic, but it could also be about a different event.).
- If you can confirm that a query is about an event distinct from the one in the query, the video is *Not Relevant*.

**Not Relevant [0]** A video is *Not Relevant* if it does not address the event described in the query. This includes videos that are about similar but distinct events—for example, a video about the *2018 Winter Olympics* when the query is about the *2022 Winter Olympics*.

**When to escalate** Escalate in the following scenarios:

- You are unable to access or play a video (before you escalate for this reason, try to open the video in another window). [Select *Video Removed*]
- The video is in the wrong language (in relation to the project you are working on)
- Given your understanding of the task and these guidelines, you are unable to determine the video's relevance in relation to one or more queries. In this scenario, please also post a message to the Teams chat describing your question. If your question is answered, you can return to the task later, update judgments as needed, and un-select the Escalate option. [Select *Escalate*]

Note that you can still provide relevance judgments for some queries, even if you escalate the task as a whole. If there are queries for which you can confidently judge relevance, please do so. If you escalate a task, leave a brief explanation in the comment box to describe your reason for escalation and any other details you think would be helpful for the next person evaluating the video.