# The Open World of Micro-Videos

Phuc Xuan Nguyen[1], Gregory Rogez[2], Charless Fowlkes[1], Deva Ramanan[3]

UC Irvine[1], INRIA[2], Carnegie Mellon University[3]

**Abstract.** Micro-videos are six-second videos popular on social media networks with several unique properties. Firstly, because of the authoring process, they contain significantly more diversity and narrative structure than existing collections of video "snippets". Secondly, because they are often captured by hand-held mobile cameras, they contain specialized viewpoints including third-person, egocentric, and self-facing views seldom seen in traditional produced video. Thirdly, due to to their continuous production and publication on social networks, aggregate micro-video content contains interesting open-world dynamics that reflects the temporal evolution of tag topics. These aspects make micro-videos an appealing well of visual data for developing large-scale models for video understanding. *We analyze a novel dataset of* **260 thousand** *micro-videos labeled with* **58 thousand** *tags.* To analyze this data, we introduce viewpoint-specific and temporally-evolving models for video understanding, defined over state-of-the-art motion and deep visual features. We conclude that our dataset opens up new research opportunities for large-scale video analysis, novel viewpoints, and open-world dynamics.

We examine an increasingly prevalent form of media known as *micro-videos*, time-constrained (typically 5-10 second) video clips commonly used on social networking sites such as Instagram and Vine. Micro-videos can be interpreted as the visual analog of a character-limited micro-blogs or "tweets" [1]. An estimated 12 million micro-videos are posted to Twitter each day. The number of micro-videos produced *surpasses the total inventory of YouTube every 3 months*. From an applied perspective, this flood of visual data is increasingly important and has unique characteristics that are not addressed by existing computer vision methodologies and benchmarks. We further argue that the microvideo format offers unique opportunities for basic research in building systems that address lifelong learning in *open-world* visual domains.

**Ease of collection/processing:** A particularly attractive aspect of micro-videos is the ease of large-scale collection, storage and processing. While large-scale datasets [2,3] have revolutionized static-image processing, the counterpart for video-based recognition does not appear to exist. One reason is that it is notoriously challenging to collect, store, and process a large diverse video collection because of resource constraints. Indeed, existing video collections often contain multiple snippets cropped from a few longer videos (to simplify the collection process) [4,5]. As we experimentally validate, in Sec. 4.2, this limits their diversity when compared to micro-video collections.

Fig. 1: Micro-videos lie in a regime between images and traditional videos, encoding semantically rich micro-narratives while remaining tractable to collect, store and process. Mobile-videographers often interact with the scene and its subjects resulting in a wide range of camera viewpoints including **egocentric** views of activities and **self-facing** shots where a single individual is both the photographer and subject. Our dataset, **MV-58k**, includes common tags about actions and objects seen in other computer vision datasets as well as specific tags such as `#noseguitar` and `#puppetman` which are video-graphic styles unique to the micro-video sharing service Vine. Distribution of videos per tag is highly skewed. At the time of submission our dataset totals more than 250,000 videos and includes at least 1000 videos for 50% of the tags.

**Micro-narratives:** An intriguing aspect of micro-videos is that they arrive "ready-for-analysis", due to the constraint that users are forced to trim content to fit within the six-second restriction. Micro-videos often contain smaller *nano-shots*, each sometimes containing a unique viewpoint, that are spliced together so as to contain a complete narrative (Fig. 5). Indeed, micro-videos are now given their own categories in established film festivals [6]. As a result of this intense degree of editing, each frame of a micro-video typically has high information content compared to frames in a longer unconstrained video. This changes strategies for automated understanding and may even eliminate the need for common video preprocessing steps like keyframe or shot selection.

**Viewpoint:** A unique technical aspect of micro-videos found on social networks is camera viewpoint (Fig. 4). Current action datasets in computer vision focus on *third-person* depictions of actions, where a person(s) performing an action or activity is framed in the view. In contrast, a significant fraction of socially-driven micro-videos include *egocentric* viewpoints, where the photographer is participating in the action. Another camera configuration is a *self-facing* viewpoint, or "selfie", where a single user is both the photographer and subject. This is particularly interesting in the study of interactions photographers and their subjects [7]. Indeed, the "selfie" is commonly recognized as a medium for embodiment and empowerment because of the unique photographer-subject interaction it afford [8]. Such diverse camera viewpoints represent new modes of video acquisition that are not typically addressed in previous work and deserve a closer look from the computer vision community.

**Open-world dynamics:** Micro-videos come labeled with hashtags that enable search and play an important role in social communication. These tags

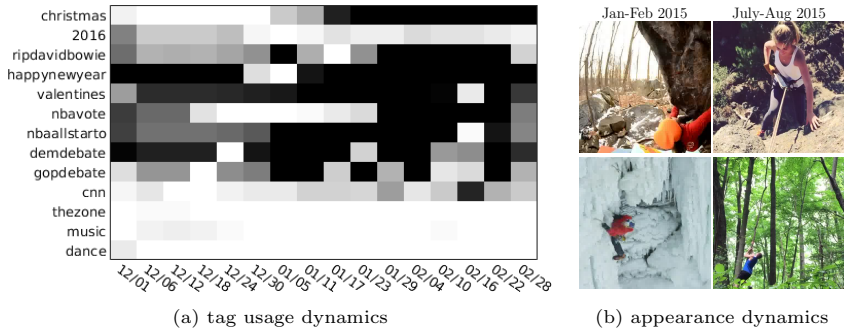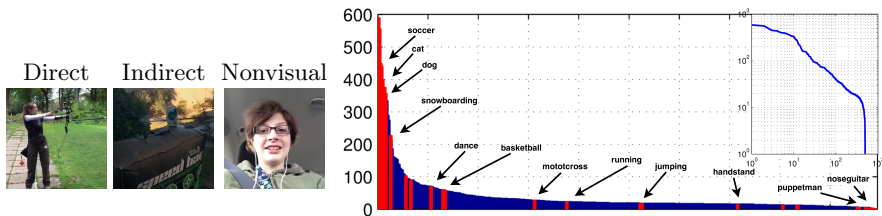(a) tag usage dynamics                    (b) appearance dynamics

Fig. 2: Unique temporal structure in microvideos. We visualize changes in tag label priors over time by plotting their popularity rank as a heatmap (brighter denotes higher popularity). Note that many tags exhibit large fluctuations; `#gopdebate` is ranked first during the weeks of the associated events, but dramatically decreases otherwise. We can also interaction between the tags; `#cnn` rises in popularity as `#demdebate` and `#gopdebate` become popular. Some tags (`#music`) exhibit relatively stable popularity. The distribution of visual appearance associated with some tags also exhibit temporal dynamics. `#climbing` shows large temporal variations over summer and winter months due to changes in scenery (desert rocks versus ice) and equipment.

provide a form of supervision, automatically labeling our diverse, multi-view, pre-trimmed dataset. In contrast to existing efforts to develop top-down ontologies for activities [9] or events [10], tags form an open-world vocabulary whose usage and semantics changes dynamically over time [11]. For example, the tag `#trump2016` did not exist until recently, while the visual meaning of the tag `#apple` expands whenever a new iPhone is released. Microvideos thus provide a unique opportunity to explore learning *in the open world*, where distributions of visual semantics follow long-tail statistics that change time over time. Taking a *bottom-up, data-driven* approach to video content semantics more accurately reflects naturally-occuring long-tail distributions that typically suppressed in hand-curated datasets. The dataset we present is thus large, has dynamic temporal variations and is continually growing is size. We analyze a snapshot as of Feb 2016 consisting of 264,327 videos with 58,243 tags. Based on a conservatice estimate of the current growth rate, the dataset will approximately include 700,000 videos and 120,000 tags by the time of ECCV.

# 1   Related Work

There is little existing work on micro-video analysis, as the medium itself is new. Redi *et al.* [12] explore the problem of finding creative micro-videos, inspired by similar studies of image quality assessment. Sano *et al.* [13] analyze the problem of detecting loop micro-videos (that are designed to be played in a continuous six-second loop). Here we focus on analyzing general properties of micro-videos, with the explicit goal of constructing a new, large-scale benchmark for temporally-evolving tag prediction.

**Viewpoint modeling:** A unique contribution of our dataset is the diversity of camera viewpoints. Existing video benchmarks for action recogni-

Direct    Indirect    Nonvisual

(a) #archery videos                    (b) Long-tail distribution of tags

Fig. 3: (a) Videos tagged with `#archery` may show direct evidence of the tagged activity, circumstantial evidence (video of a can which is suddenly pierced by an arrow), or non-visual evidence (a girl talking about an archery competition). For diagnostic purposes we construct a hand-curated dataset (MV-40) containing only videos with direct visual evidence for a subset of 40 tags selected to span both common and rare tags, as shown in the long-tailed distribution (b). Inset shows the same distribution on a log-log scale.

tion have focused on third-person viewpoints (e.g., HMDB [4], UCF101 [14], Hollywood-2 [15], UT-Interaction [16] and Olympic Sports [17]). Wearable cameras such as Google Glass and GoPro have spurred interest in analyzing egocentric views [18,19,20,21]. Compared to existing action and egocentric datasets, our user-generated micro-videos contain a wider variety of categories (tags) and viewpoints (e.g., self-facing) with richer narrative content, even in a single clip. To understand and highlight these differences we train mixture-of-viewpoint models that specifically target viewpoint variations in dynamic micro-videos (Sec. 3) and carry out an extensive comparison with HMDB [4] (Sec. 4).

**Closed vs open vocabularies:** Traditionally, video datasets in computer vision have been labeled with a fixed ontology of activities or events [9,10,22]. An alternative perspective (popular in the multimedia community) is to formulate the problem as a multi-label tag or concept prediction task [23,24,25,26,27]. Our dataset falls into this later camp. In terms of size and diversity, the most relevant prior work appears to be Sports-1M [22] which contains 1M videos in 487 categories, and EventNet [10], which contains 95K videos labeled with 5K concepts. Our dataset already includes 2X more videos and 10X more concept tags. Unlike other video datasets, our data also includes timestamps which allow us to study temporally-varying semantics, a relatively unexplored concept in vision, with the notable exception of [28]. Importantly, tag frequency distributions are highly imbalanced, following a natural long-tail distribution (Fig. 3). While highly imbalanced class distributions are somewhat uncommon in current vision datasets, they appear to be a fundamental aspect of life-long learning in the open-world [29]. With the advent of deep architectures that appear capable of transferring knowledge across imbalanced classes [30], we think the time is right to (re)consider learning in the open-world!
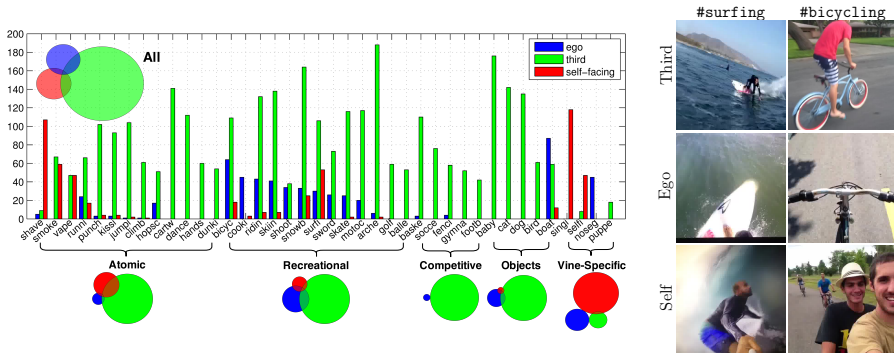
Fig. 4: The unique viewpoints of microvideos. (**a**) shows the distribution of camera viewpoints annotated in MV-40. Most tags tend to be associated with third-person viewpoints. Some atomic actions such as #shave, #smoke, and #vape (smoking with an electronic cigarette) are captured with self-facing views. Tags for recreational sports often include egocentric views (where the photographer is engaged in the action), but this is less likely for competitive sports (since it may be difficult to hold the camera). Many Vine-specific tags make use of self-facing viewpoints. From the Venn diagram embedded in the bar graph, we can see that there's a significant amount of videos that contain more than one viewpoint and some has all 3 viewpoints. We show examples of frames with unique viewpoints in (**b**). The biking frame could be labeled as both *ego* and *third* because the photographer and subjects are engaged in a "social" activity. We also show a *self-facing* viewpoint made possible through a specialized camera mount.

## 2  Dataset

In this section, we describe our (continually-running) data-collection process and analyze the statistics of micro-video tags, shots and views that make our dataset distinct from existing video benchmarks.

**Streaming dataset collection:** We collect a stream of Vine videos by daily querying of Vine's API [31]. To ensure a diverse stream, we query the 300 most popular and 300 most recent videos across multiple community-curated channels (Comedy, Sports, Musics). We typically obtain  6000 videos daily. Each video is associated with a collection of hashtags that are added to our open vocabulary. On average, a video contains 1.56 hashtags, but this statistic is skewed by the fact that nearly half the videos do not contain any tags (56%).This indirectly motivates one practical application of our dataset - automatic prediction of hashtags. We attempted to automatically merge similar tags through linguistic normalization, but found this merging did little to change our label space (perhaps because users have an incentive to used normalized hashtags that are already searchable). We visualize our overall distribution of tags in Fig. 1. We use **MV-58k** to refer to a snapshot of this datastream collected during the period Dec-2015 to Feb-2016. We provide additional statistics of this data collection and projections in the supplemental materials.

**Curation:** To examine the amount of "noise" in the tag stream and perform diagnostic comparisons to existing activity datasets, we manually curated a subset of 40 tags and their associated 4000 videos. We selected 40 representative tags
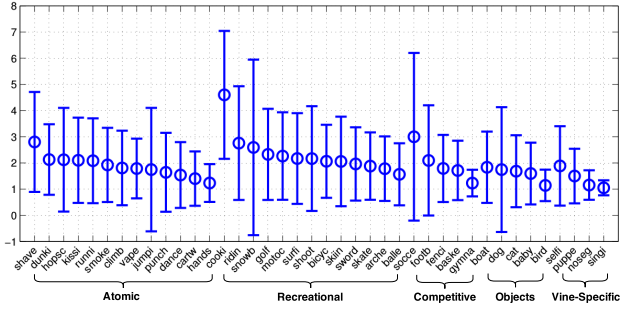
Fig. 5: Unique shot structure in microvideos. We plot the mean and variance of the number of shots in a video. Over 34% of our videos contain more than a single shot. A small percentage(3.75%) contain 6 or more shots, implying that many shots are less than 1 second in duration. In terms of per-tag statistics, `#cook`ing contains a large number of spliced shots, necessary to temporally compress and summarize such a long duration activity. `#snowboard`ing and `# soccer` have a large variance because they contain some compilation videos composed of many shots.

that span both common actions as well rare tags "in-the-tail". We "clean" this dataset by merging synonymous tags (e.g., `#horseback` and `#horsebackriding`), removing mistagged (spam) videos, and removing videos in which there was only circumstantial visual evidence for the tagged activity (see Fig. 3 for an example). We added additional annotations to each video including viewpoint and a dominant (tag) category. The latter allows us to recast tag prediction as a K-way classification problem, simplifying our diagnostic analysis. We refer to this curated dataset as MV-40, and contrast this with (the fixed snapshot of) our uncurated, open-world dataset MV-58K. We organize MV-40 into broad categories of atomic actions, recreational activities, competitive sports, objects, and vine-specific. We visualize our two-level taxonomy and provide visual examples in Fig. 1.

**Long-tail distributions:** Our collection process reveals a salient property of open-world microvideos; they follow long-tail distributions of tags. This significantly complicates learning because there will be some tags for which we have little training data. While traditionally a notorious challenge for machine learning, our analysis suggests that hierarchical feature learning (with CNNs) can learn to share, or transfer knowledge from the data-rich tags to the data-sparse tags (i.e., *one-shot learning*). For example, even if we have few examples of `#dunking`, mid-level features learned for a data-rich tag such as `#basketball` may still be useful for the former class.

**Temporal dynamics:** Our open-world dataset collection has another notable property - both the frequency of tag usage and the visual appearance semantics associated with a given tag evolve over time. In probabilistic terms, we can interpret such dynamics and changes in the *prior* of labels (Fig. 2) and the *likelihood* of image features conditioned on the tag label (Fig. 10). This suggests developing approaches for continually retraining models as new data becomes available, an idea we explore in Sec.4. An extreme case arises when a new tag first appears in the data stream there are no training examples available (e.g.,

`#trump2016` did not exist until recently). In our current experiments, we simply fail to predict such tags at test time. However, we point out that the temporal appearance of new tags provides a compelling natural example of *zero-shot learning* where side information about the semantic relation between tags could readily be exploited.

**Shot statistics:** Micro-videos have a unique shot structure due to the capturing mechanics and limited time constraint imposed by the app. A Vine or Instagram user creates videos by holding a button to capture and releasing the button to pause at any time. He/she can later resume capturing again until the content time limit is exceeded (6 seconds for Vine). This interface allows and even encourages splicing together of many shots within a single time-limited video. We plot a distribution shot length and frequency for various tags in Fig. 5. Some tags, such as `#cooking` tend to consistently involve a large number of shots. Indeed, we find 3.75% of videos have more than 6 shots, implying many shots are less than a second in length.

**Viewpoint:** To analyze the effect of viewpoint on recognition accuracy, we manually annotate MV-40 with the viewpoint of each video as egocentric, third-person, or self-facing. In some cases, the viewpoint changed between shots in the video in which case we record all the viewpoints present, as well as the dominant one. With this annotation we can treat viewpoint prediction as either a multi-label attribute prediction problem (where multiple viewpoints can be present in video) or a multi-class problem (where the goal is to predict the dominant viewpoint). We report the per-tag viewpoint statistics in Fig. 4. Interestingly, even at the shot-level, we sometimes find ambiguities in viewpoint. Some social activities (such as `#bicycling`) involve both the photographer and subjects in view, suggesting a simultaneous egocentric and third-person viewpoint.

## 3   Methods

In this section, we describe several baseline models for tag prediction and viewpoint classification. As the focus of our work is not on video feature extraction, we make use of standard feature sets. Recent results from the THUMOS evaluation benchmark [32] suggest that CNN spatial features (VGG [33]) combined with temporal motion features (IDT [34]) make for a reasonable video descriptor. We briefly outline our feature descriptor and associated classification engines here.

**Appearance features:** Recent work has shown that Convolutional Neural Networks (CNN) produce quite effective visual descriptors for recognition tasks [35,33] including video analysis [36,37,22,38]. We experimented with many open-source implementations of CNN architectures for video-feature extractors, and found that while many were effective, some placed significant demands on processing time and descriptor storage. We refer the reader to our supplementary material for a detailed description of our diagnostic experiments. We found a good tradeoff in speed, storage, and accuracy with the following simple pipeline: given a video, (1) run off-the-shelf VGG-16 models [33] on 15 equally-spaced

frames from a video, (2) for each frame, extract (6144-dimensional) multi-scale features across multiple layers [39], and (4) max-pool the resulting features across the 15 frames [36]. When compared to standard single-frame CNN feature extraction (that extract 4096-dimensional features [33]), our final video pipeline was only 15x slower and 2x larger in storage costs.

**Motion features:** Recent state-of-the-art results on video datasets have made use of trajectory-based motion features. We include such features in our analysis, focusing on Improved Dense Track (IDT) [34]. This method is based on a "interest-track" framework (rather than a space-time interest-point approach) in which short-term tracks are found from tracking interest points across frames. One then extracts various features aligned to these temporal tracks, including oriented gradient histograms and optical flow. We quantize these descriptors based on established guidelines for constructing a codebook, using $K = 4000$ codebook entries found with K-means.

**Cue-combination:** We train tag classifiers that aggregate features by combining multiple kernels. We experimented extensively with various feature encodings and kernel combinations before settling on the following strategy. We compute the similarity between two video clips $x_i$ and $x_j$ by averaging their appearance and motion-feature similarities:

$$K(x_i, x_j) = \frac{1}{2}(K_m(x_i, x_j) + K_a(x_i, x_j))$$

We define the motion similarity with a $\chi^2$-RBF kernel:

$$K_m(x_i, x_j) = exp\left(-\frac{1}{L}\sum_{c=1}^{L}\frac{d(x_i^c, x_j^c)}{A^c}\right), \quad x_i^c, x_j^c \in \mathcal{R}^{4000}$$

where $A_c$ is the average $\chi^2$ distance between all videos in the training data, L is the number of motion channels (Traj, HOG, HOF, MBHx, MBHy), and $d(x_i^c, x_j^c)$ is the $\chi^2$ distance between $x_i$ and $x_j$ with respect to the $c$-th channel. We measure the appearance similarity between two clips using a linear kernel: $K_a(x_i, x_j) = \sum_{f=1}^{N}(x_i^f \cdot x_j^f)$ summed over the $N = 15$ static feature channels extracted by the CNN. Given a training vocabulary of $K$ tags, we train $K$ binary (one-vs-all) kernelized SVMs using the LIBSVM package [40]. Finally, we calibrate each predictor using Platt scaling [41].

**Viewpoint mixtures:** Different viewpoints of the same tag can have significant differences in video content, as shown in Fig. 1. For example, a third-person and egocentric `#bicycling` video contain very different motions and appearances. To analyze such variations in the MV-40 diagnostic dataset, we train viewpoint-specific models for each tag. The final confidence associated with a tag prediction is the maximum score across the three viewpoint-specific models.

**Temporally-adapted models:** To explore the temporal evolution of tag semantics and videos, we evaluate models trained with videos sampled over different temporal windows. Consider the task of predicting the label $y$ for a video $x$ collected at time $t_0$ using with models trained on a stream of timestamped

training videos indexed by time. We consider models trained on three different subsets of training data:

$$\{x_t, y_t : \forall t\} \rightarrow \text{non-causal model} \tag{1}$$

$$\{x_t, y_t : t < t_0\} \rightarrow \text{causal model} \tag{2}$$

$$\{x_t, y_t : t_0 - \Delta < t < t_0\} \rightarrow \text{adaptive model} \tag{3}$$

where $\Delta$ is a specified window size. The above approach can be simplified by making some assumptions about the nature of temporal variation. For example, if we assume that the popularity of tags changes over time but their appearance models do not, one can model efficiently model temporal dynamics with a *statistical prior shift* [42]. Intuitively, dynamics can be captured with a fixed set of posterior class predictions that are reweighted by dynamically-varying tag priors. Unfortunately, this requires access to tag priors on test data from the future, which violates causality. Instead, we assume that tag priors vary smoothly over time, and simply use a weighted estimate of recent tags' popularity. We found that the simple approach of applying temporally-weighted Platt rescaling (using a weighted dataset where recent videos are given more importance) outperformed an explicit prior model.

## 4   Experiments

In this section, we present an extensive set of experiments on our dataset and refer the reader to the supplementary materials for additional tables and figures. We focus on three sets of experiments: a diagnostic evaluation of features and viewpoints on our curated dataset (MV-40), its relation to popular benchmarks such as HMDB [4], and analysis of the open-world dynamics of MV-58K.

### 4.1   Diagnostics

First, we analyze various aspects of our dataset and recognition pipeline, focusing on the curated and annotated MV-40 subset.

**Feature comparisons:** We begin by comparing the performance of various combinations of our features in Fig. 6-(a). We observe CNN features outperform IDT in most category groups except 'Atomic'. Trajectory-based motion and appearance-based deep features are particularly effective when combined, indicating that they capture complementary cues. Supplementary materials include the class confusion matrix over all 40 tags for the IDT+CNN feature combination.

**View-specific mixtures:** We next evaluate the performance of view-specific tag classifiers and compare to results with a single classifier per tag. To ensure sufficient training data for each mixture component, we only train a mixture for a specific view if there are more than 20 videos in that viewpoint (11 classes out of 40 satisfy this criteria in MV-40). When training, one can treat clips from the same tag but different viewpoints as positive training examples (**Pos**), negative

| | IDT | CNN | IDT+CNN |
|---|---|---|---|
| All | 53.90 | 62.13 | 68.82 |
| Atomic | 47.72 | 46.58 | 57.19 |
| Recreational | 60.81 | 69.09 | 76.06 |
| Competitive | 49.86 | 60.29 | 64.64 |
| Objects | 46.14 | 70.37 | 73.47 |
| Vine-Specific | 63.25 | 65.81 | 68.38 |
| Ego | 63.25 | 73.15 | 79.49 |
| Third | 51.45 | 60.00 | 67.14 |
| Self | 59.27 | 61.82 | 68.40 |

| Feature | Pos | Neu | Neg |
|---|---|---|---|
| IDT | 62.65 | 64.98 | 64.34 |
| CNN | 58.28 | 60.82 | 60.47 |
| IDT+CNN | 68.57 | 70.61 | 70.40 |

(a) Feature performance        (b) View-specific mixtures        (c) Viewpoint confusion

Fig. 6: (a) Performance of different feature sets on MV-40 broken down by tag type and viewpoint. (b) plots the perfomance of view-point specific tag models, exploring different choices of positive and negative data. (c) visualizes the 3-way class confusion matrix for viewpoint prediction. See text for more details.

training examples (**Neg**), or such clips can be treated as neutral and ignored (**Neu**). Fig 6-(b) summarizes of the performance. When averaged over all tags, the performance increase from view-specific mixtures is rather negligible ( 0.6% for our combined features). We also evaluate only those 11 classes for which additional views were trained. We see a small but definite improvement of 2.1%. Ignoring clips from other viewpoints (**Neu**) slightly increases performance.

**Viewpoint prediction:** We also investigate the task of viewpoint prediction: what is the viewpoint of a test video? Fig 6-(c) summarizes viewpoint confusions. Egocentric views are often confused with third-person and the accuracies for egocentric drops significantly. This is consistent with Fig. 4, which suggests many recreational activities involve both the photographer and subjects involved in the action. If we score viewpoint prediction as a multilabel problem (where each video could be labeled with more than one viewpoint), accuracy for `ego`, `third`, and `self-facing` jump to 92%, 90%, and 93%. This suggests the presence of any given viewpoint can be accurately predicted.

## 4.2   Comparison to existing benchmarks

We perform an extensive comparison of our data with a popular action recognition benchmark, HMDB [4]. We use a subset of 15 tags (**vine-15**) that overlap with HMDB categories and evaluate the IDT+CNN based predictor. Overall, the average accuracy on **vine-15** is lower than HMDB (65.99% vs. 71.27%), suggesting that our data is more challenging. Torralba and Efros [43] suggest that the 'performance drop' provides a way to quantify how biased or general a dataset may be. The drop for models trained on Vine data is 13%, while the drop for models trained with HMDB is 26.75%.

**Viewpoint:** One might hypothesize that since HMDB contains mostly third-person views, it won't generalize to the other viewpoint in our data. To test this, we extract a smaller subset **vine-3rd** containing only third-person viewpoints. Fig 7-(g) shows that **vine-3rd** is more similar to HMDB, as models trained on HMDB data perform better on **vine-3** than **vine-15**. However, performance drop for models trained on **vine-3rd** is still significantly smaller than those

trained on HMDB (12.88% vs. 24.4%). This suggests that even accounting for viewpoint, our videos still generalize better than HMDB.

**Temporally iconic videos:** The concept of "iconic views" in object recognition refers to "easy" images with a clear and distinctive depiction of an object, often close cropped or in an uncluttered setting without occlusion. We apply this notion to video, defining a *temporally iconic depiction* of a tag as one where temporal clutter has been removed by trimming down the video clip to focus on the core action. Temporal cropping is relevant even in micro-videos, which often contain additional frames and shots surrounding those described by the tag. We manually segment each video in **vine-15** to derive an iconic version **vine-icon** and use **vine-3rd-icon** to denote the third-person subset. Fig. 7 reveals that **vine-icon** and **vine-3rd-icon** are slightly easier than vine, and more similar to HMDB (following our previous analysis). However, the cross-dataset performance drop for models trained on **vine-icon** and **vine-3rd-icon** (14.06% and 12.86% respectively) are still significantly smaller than the drop by models trained **HMDB** data (25.35% and 21.89%). *To summarize, even though third-person and temporal iconic-ness accounts for much of the difference between HMDB and our dataset, our micro-videos can still generalize better.*

**Qualitative differences:** To better understand the differences, we visualize example videos in Fig. 7. #swordfighting in **vine-15** tends to involve people playfully jousting in social everyday scenes, while HMDB clips tend to involve formally sparring or staged/scripted fights. This difference helps explain the why HMDB's #swordfights don't generalize to Vine and is consistent with the notion that models trained from purely iconic images of objects often generalize poorly [43]. In the same figure, we also plot sample frames from videos with a particular tag. *Perhaps surprisingly, many distinct clips in the dataset actually come from the same longer video.* This reduces the amount of diversity in the dataset, and reinforces one of our motivations: it *is* surprisingly hard to collect diverse video clips. This phenomena is not limited to HMDB, and also appears in other benchmarks such as ImageNet Video Challenge [5]. Our micro-video dataset, however, is inherently *diverse by construction* due to its dynamic and pre-trimmed nature.

## 4.3   Open-World Dynamics

In this section, we analyze properties of our open-world MV-58K dataset. One immediate issue is that open-world tags are more naturally treated as multi-label tasks (since videos can be naturally labeled with many tags). We adopt two scoring criteria from the literature on multi-label classification [44]. $mAP_T$ treats each tag as a individual binary prediction problem, and computes the average precision (AP) of each prediction task, returning the mean AP over all tags. This approach equally weights popular and infrequent tags, and is analogous with standard mAP measures used in object detection. Alternative, $mAP_I$ first computes an image-specific AP by ranking all tag predictions specific to that image, and then returns the mean AP over all images. This latter scheme places more importance on frequent tags, and also requires that confidences across

(a) Vine swordfight



(b) HMDB swordfight



(c) Vine



(d) HMDB

| Train | Test | avg |
|---|---|---|
| vine-15 | vine-15 | 65.99 |
| vine-15 | hmdb | 52.86 |
| vine-3rd | vine-3rd | 65.82 |
| vine-3rd | hmdb | 52.94 |

(e)

| Train | Test | avg |
|---|---|---|
| vine-icon | vine-icon | 68.03 |
| vine-icon | hmdb | 53.97 |
| vine-3rd-icon | vine-3rd-icon | 68.65 |
| vine-3rd-icon | hmdb | 55.79 |

(f)

| Train | Test | avg |
|---|---|---|
| hmdb | hmdb | 71.27 |
| hmdb | vine-15 | 44.52 |
| hmdb | vine-3rd | 46.87 |
| hmdb | vine-icon | 45.92 |
| hmdb | vine-3rd-icon | 49.38 |

(g)

Fig. 7: Qualitative difference between HMDB and Vine. The **(top)** row shows frames from the `#swordfight` videos for **(a)** Vine and **(b)** HMDB. Micro-videos often contain people playfully jousting in everyday scenes, while HMDB contains more formal, staged fights and sparring events. **(Bottom)** Snapshots of `#cartwheel` videos for **(c)** Vine and **(d)** HMDB. Multiple clips from HMDB can come from one source video, but Vine videos are naturally collected from different sources. We evaluate various cross-dataset experiments between subsets of our dataset and HMDB in (e) - (g).

different tag predictors be calibrated. By default, we use $mAP_I$ unless otherwise specified.

**Curated vs raw data:** When comparing models trained on curated versus open-world data (Fig. 8), we see that additional data always helps, but given a fixed amount of training data, curated data performs better. Curated data is designed to mimic existing recognition datasets which use concise notions of visual categories, while the raw dataset captures problem of tag-prediction "in the wild". For example, videos of people talking about `#archery` are removed in MV-40, but exist in MV-58k (Fig. 3). When comparing accuracy on tag-prediction, we find that a clean dataset replicates the accuracy of a 2X larger raw dataset.

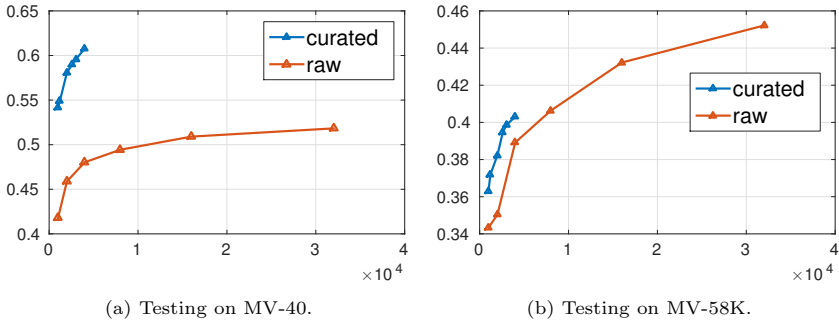(a) Testing on MV-40.                    (b) Testing on MV-58K.

Fig. 8: To evaluate the effect of data curation, we compare the accuracy of models trained and tested on the curated MV-40 versus uncurated data (obtained by targeted querying from the open-world for additional videos with the 40 given tags). The curated testset (a) more closely resembles existing recognition datasets while (b) is representative of tag prediction "in the wild". More training data improves action recognition (MV-40), but curated data is significantly more effective than a large, raw training set. For tag prediction (b), a modest amount of additional raw training data ($< 2X$) rivals the accuracy of a manually curated training set.

**Closed vs open vocabularies:** We now move beyond our 40 selected tags to evaluation of an open-world vocabulary. Because the amount of training data variables considerably per tag, we plot test accuracy as a function of training-set size in Fig. 9. We tend to see different performance regimes. "Easy" tags perform well even with little training data, likely due to a characteristic appearance that is easy to learn from little training data (`#cavs, #warriors`). "Challenging" tags appear to contain appearance variation, but are learnable with additional data (`#dogs, #soccer`). "Unlearnable" tags remain near-zero AP even given lots of training data (`#revine, #lol`). We posit that these can be treated as stopwords that fail to capture much semantic meaning of the video.

**Temporal dynamics:** We now examine the time-varying properties of micro-video content and tags (Fig. 10). We refer the reader to the caption for a detailed analysis, but we find that causal prediction (where one only has access to data from the past) is much harder than the non-causal counterpart, suggesting that micro-videos are not "iid" over time. Benchmarks that use iid resampling may over-estimate real-world performance on streaming data. We find that much of this temporal variation can be explained by fluctuations of tag popularity, whose degree of variation can vary dramatically across classes.

**Conclusion:** We introduce a open-world dataset of micro-videos, which lie in a regime between single images and typical videos, allowing for easy capture, storage, and processing. They contain micro-narratives captured from viewpoints typically not studied in computer vision. Because they are naturally diverse, pre-trimmed, and user-annotated, they can be used a live testbed for open-world evaluation of video understanding systems. We conclude with an intriguing thought: rather than distributing a fixed benchmark dataset (which historically leads to eventual overfitting [45]), our analysis suggests that we can instead
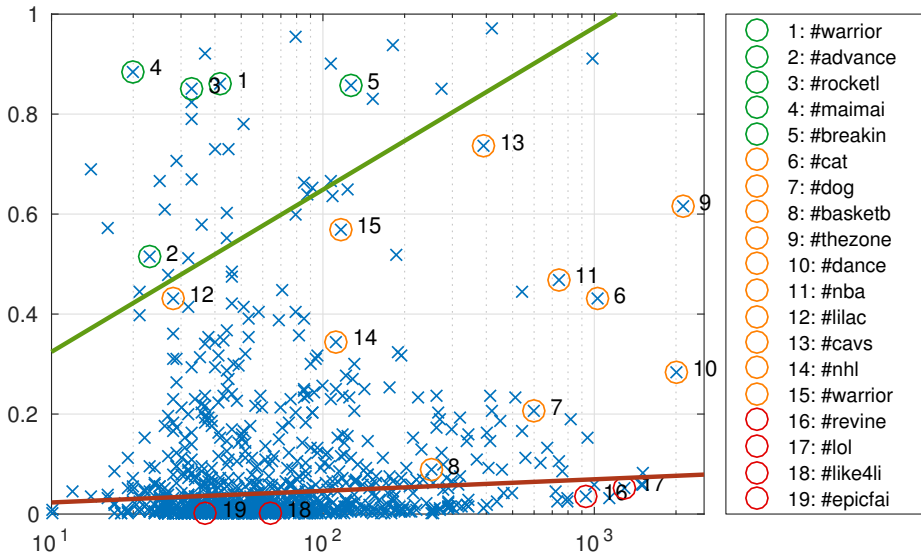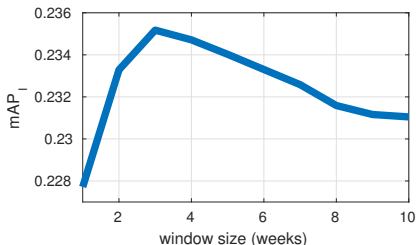
Fig. 9: Open-world tag prediction. The scatter plot shows per-tag APs vs (log) number of training examples. We rank the "learnability" of a tag by the ratio of its AP to (log) number of training examples, and draw lines to loosely denote regions of easy, challenging, and unlearnable tags. Unlearnable tags appear to correspond to "stopwords" such as `#revine,lol` that do not capture video content. The per-tag $mAP_T$ is 0.05.

| Category | Non-causal | Causal | Adaptive |
|---|---|---|---|
| 40 categories | 44.55 | 38.05 | 43.60 |
| `#climbing` | 6.87 | 5.27 | 21.16 |

(a)



(b)

Fig. 10: Temporal dynamics. We plot performance of different temporal training strategies in (**a**). Non-causal models perform the best but may be impractical because they must be trained on future videos. Adaptive models trained on recent videos perform better than a fixed training set because tag topics tend to temporally evolve. In (**b**), we analyze a hybrid causal-adaptive strategy that learns causal classifiers with adaptive Platt calibration (calibrated on the past $w$ weeks). With the window size at $w = 3$, the performance improves to 23.5% (from 21% with no calibration or calibrated on training data). If we calibrate on the ideal test distribution in the future, performance jumps significantly to 28%, indicating that temporal models that can predict *future* popularity of tags would significantly improve accuracy..

distribute a benchmark script that evaluates models on live open-world microvideos. We think the time is right to consider video recognition out in-the-open!

# References

1. Kwak, H., Lee, C., Park, H., Moon, S.: What is twitter, a social network or a news media? In: WWW Conference, ACM (2010) 591–600
2. Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes challenge: A retrospective. IJCV **111**(1) (2014) 98–136
3. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR 2009, IEEE (2009) 248–255
4. Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., Serre, T.: HMDB: a large video database for human motion recognition. In: ICCV, 2011
5. Liu, W., Berg, A.: Imagenet video object detection. `http://image-net.org/challenges/LSVRC/2015/index`
6. Kelly, P.: Slow media creation and the rise of instagram. Mobile Media Making in an Age of Smartphones (2014) 129
7. Marien, M.W.: Photography: A cultural history. Laurence King London (2002)
8. Jones, A.: The eternal return: Self-portrait photography as a technology of embodiment. Signs **40**(1) (2014)
9. Fabian Caba Heilbron, Victor Escorcia, B.G., Niebles, J.C.: Activitynet: A large-scale video benchmark for human activity understanding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2015) 961–970
10. Ye, G., Li, Y., Xu, H., Liu, D., Chang, S.F.: Eventnet: A large scale structured concept library for complex event detection in video. In: ACM MM. (2015)
11. Cunha, E., Magno, G., Comarela, G., Almeida, V., Gonçalves, M.A., Benevenuto, F.: Analyzing the dynamic evolution of hashtags on twitter: a language-based approach. In: Proceedings of the Workshop on Languages in Social Media, Association for Computational Linguistics (2011) 58–65
12. Redi, M., OHare, N., Schifanella, R., Trevisiol, M., Jaimes, A.: 6 seconds of sound and vision: Creativity in micro-videos. In: CVPR, 2014
13. Sano, S., Yamasaki, T., Aizawa, K.: Degree of loop assessment in microvideo. In: ICIP 2014, IEEE (2014) 5182–5186
14. Soomro, K., Zamir, A.R., Shah, M.: Ucf101: A dataset of 101 human actions classes from videos in the wild. arXiv preprint arXiv:1212.0402 (2012)
15. Marszalek, M., Laptev, I., Schmid, C.: Actions in context. In: CVPR, 2009
16. Ryoo, M., Aggarwal, J.: Ut-interaction dataset, icpr contest on semantic description of human activities (sdha) (2010)
17. Niebles, J.C., Chen, C.W., Fei-Fei, L.: Modeling temporal structure of decomposable motion segments for activity classification. In: ECCV, 2010
18. Fathi, A., Farhadi, A., Rehg, J.M.: Understanding egocentric activities. In: ICCV, 2011
19. Kitani, K.M., Okabe, T., Sato, Y., Sugimoto, A.: Fast unsupervised ego-action learning for first-person sports videos. In: CVPR 2011, IEEE (2011) 3241–3248
20. Pirsiavash, H., Ramanan, D.: Detecting activities of daily living in first-person camera views. In: CVPR 2012, IEEE (2012) 2847–2854
21. Lee, Y.J., Grauman, K.: Predicting important objects for egocentric video summarization. IJCV (2014) 1–18
22. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L.: Large-scale video classification with convolutional neural networks. In: CVPR, IEEE (2014) 1725–1732

23. Vahdat, A., Mori, G.: Handling uncertain tags in visual recognition. In: ICCV 2013, IEEE (2013) 737–744
24. Yang, W., Toderici, G.: Discriminative tag learning on youtube videos with latent sub-tags. In: CVPR 2011, IEEE (2011) 3217–3224
25. Aradhye, H., Toderici, G., Yagnik, J.: Video2text: Learning to annotate video content. In: ICDMW'09, IEEE (2009) 144–151
26. Toderici, G., Aradhye, H., Pasca, M., Sbaiz, L., Yagnik, J.: Finding meaning on youtube: Tag recommendation and category discovery. In: CVPR 2010, IEEE (2010) 3447–3454
27. Ulges, A., Koch, M., Borth, D., Breuel, T.M.: Tubetagger-youtube-based concept detection. In: ICDMW'09, IEEE (2009) 190–195
28. Kim, G., Xing, E.P., Torralba, A.: Modeling and analysis of dynamic behaviors of web image collections. In: Computer Vision–ECCV 2010. Springer (2010) 85–98
29. Chen, X., Shrivastava, A., Gupta, A.: Neil: Extracting visual knowledge from web data. In: Proceedings of the IEEE International Conference on Computer Vision. (2013) 1409–1416
30. Bengio, Y., Delalleau, O.: On the expressive power of deep architectures. In: Algorithmic Learning Theory, Springer (2011) 18–36
31. : https://github.com/starlock/vino/wiki/api-reference
32. Gorban, A., Idrees, H., Jiang, Y., Zamir, A.R., Laptev, I., Shah, M., Sukthankar, R.: Thumos challenge: Action recognition with a large number of classes (2015)
33. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. ICLR (2015)
34. Wang, H., Schmid, C.: Action recognition with improved trajectories. In: ICCV, 2013
35. Razavian, A.S., Azizpour, H., Sullivan, J., Carlsson, S.: Cnn features off-the-shelf: an astounding baseline for recognition. In: CVPRW 2014, IEEE (2014) 512–519
36. Ryoo, M.S., Rothrock, B., Matthies, L.: Pooled motion features for first-person videos. In: CVPR, IEEE (2015)
37. Simonyan, K., Zisserman, A.: Two-stream convolutional networks for action recognition in videos. In: Advances in Neural Information Processing Systems. (2014) 568–576
38. Z. Xu, Y. Yang, A.H.: A discriminative cnn video representation for event detection. In: CVPR, IEEE (2015) 1–8
39. Yang, S., Ramanan, D.: Multi-scale recognition with dag-cnns. arXiv preprint arXiv:1505.05232 (2015)
40. Chang, C.C., Lin, C.J.: Libsvm: a library for support vector machines. ACM Transactions on Intelligent Systems and Technology (TIST) **2**(3) (2011) 27
41. Platt, J., et al.: Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. Advances in large margin classifiers **10**(3) (1999) 61–74
42. Saerens, M., Latinne, P., Decaestecker, C.: Adjusting the outputs of a classifier to new a priori probabilities: a simple procedure. Neural computation **14**(1) (2002) 21–41
43. Torralba, A., Efros, A.A.: Unbiased look at dataset bias. In: Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on, IEEE (2011) 1521–1528
44. Makadia, A., Pavlovic, V., Kumar, S.: A new baseline for image annotation. In: Computer Vision–ECCV 2008. Springer (2008) 316–329
45. Ponce, J., Berg, T.L., Everingham, M., Forsyth, D.A., Hebert, M., Lazebnik, S., Marszalek, M., Schmid, C., Russell, B.C., Torralba, A., et al.: Dataset issues in

object recognition. In: Toward category-level object recognition. Springer (2006) 29–48