

Bi-Calibration Networks for Weakly-Supervised Video Representation Learning

Fuchen Long, Ting Yao, *Senior Member, IEEE*, Zhaofan Qiu, Xinmei Tian, *Member, IEEE*, Jiebo Luo, *Fellow, IEEE*, and Tao Mei, *Fellow, IEEE*

Abstract—The leverage of large volumes of web videos paired with the searched queries or surrounding texts (e.g., title) offers an economic and extensible alternative to supervised video representation learning. Nevertheless, modeling such weakly visual-textual connection is not trivial due to query polysemy (i.e., many possible meanings for a query) and text isomorphism (i.e., same syntactic structure of different text). In this paper, we introduce a new design of mutual calibration between query and text to boost weakly-supervised video representation learning. Specifically, we present Bi-Calibration Networks (BCN) that novelly couples two calibrations to learn the amendment from text to query and vice versa. Technically, BCN executes clustering on all the titles of the videos searched by an identical query and takes the centroid of each cluster as a text prototype. The query vocabulary is built directly on query words. The video-to-text/video-to-query projections over text prototypes/query vocabulary then start the text-to-query or query-to-text calibration to estimate the amendment to query or text. We also devise a selection scheme to balance the two corrections. Two large-scale web video datasets paired with query and title for each video are newly collected for weakly-supervised video representation learning, which are named as YOVO-3M and YOVO-10M, respectively. The video features of BCN learnt on 3M web videos obtain superior results under linear model protocol on downstream tasks. More remarkably, BCN trained on the larger set of 10M web videos with further fine-tuning leads to 1.6%, and 1.8% gains in top-1 accuracy on Kinetics-400, and Something-Something V2 datasets over the state-of-the-art TDN, and ACTION-Net methods with ImageNet pre-training. Source code and datasets are available at <https://github.com/FuchenUSTC/BCN>.

Index Terms—Video Representation Learning, Weakly-supervised Learning, Action Recognition.

1 INTRODUCTION

WITH the rise of deep learning technologies, there has been a steady momentum of breakthroughs on video representation learning [1], [2], [3], [4], [5], [6], [7], [8], [9]. The achievements rely heavily on the requirement to have large amount of labeled data for fully-supervised training. In practice, acquiring the annotations of videos is very expensive and time-consuming. Therefore, recent research [10], [11] study the alternative regime of web data, which is largely available and freely accessible by search engines, for weakly-supervised learning. These approaches usually treat the weakly visual-textual connection as a reliable signal and directly maximize the similarity between them or alleviate the challenge of noise (e.g., incorrect or irrelevant text) in web data for training, but seldom explore the inherent property of videos or texts. For example, is it reasonable that the representations of all the searched videos from an identical query cluster around a single prototype? or if the surrounding texts of two videos are close in representation space, should the representations of the two videos also be in proximity?

- Ting Yao is the corresponding author.
- Fuchen Long, Ting Yao, Zhaofan Qiu, and Tao Mei are with JD Explore Academy, Beijing, China. Xinmei Tian is with University of Science and Technology of China, Hefei, China. Jiebo Luo is with University of Rochester. (e-mail: longfc.ustc@gmail.com; tingyao.ustc@gmail.com; zhaofanqiu@gmail.com; tmei@jd.com; xinmei@ustc.edu.cn; jluc@cs.rochester.edu).
- This work was supported by the National Key R&D Program of China under Grant No. 2020AAA0108600.

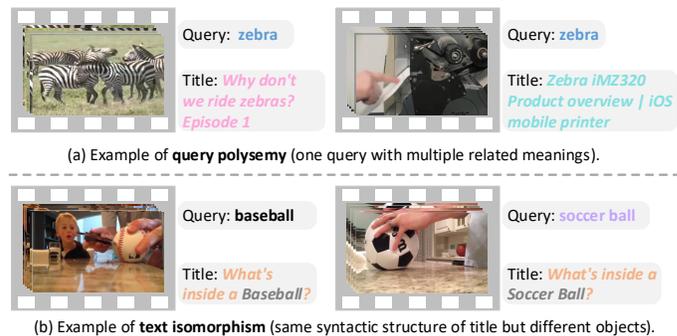


Fig. 1: The intuitive examples of (a) “query polysemy” and (b) “text isomorphism” of web videos.

In order to answer the two questions, let’s look at the two examples illustrated in Figure 1. The upper one shows two returned videos when searching for the query of “zebra” in a text-based video search engine. The query “zebra” de facto contains two types of search intention: an animal or a brand of printer, and the videos in response to the two meanings are quite distinct in visual appearance. The phenomenon is known as “query polysemy.” In this case, taking the two families of videos as one class will inevitably mislead video representation learning through either classification or cross-view embedding. We propose to mitigate the issue by leveraging the calibration from surrounding text (i.e., title in this work) of the videos. The titles of “why don’t we ride zebras? Episode 1” and “Zebra iMZ320 Product overview | iOS mobile printer” provide rich information about the video content and are

more descriptive. The videos are naturally grouped into two clusters based on valuable supervision of titles and each cluster corresponds to one semantic meaning, potentially making video representation learning more discriminative. In contrast, robust learning of video representation also necessitates using the query to regulate the text. Taking the lower case as an example, the titles of “What’s inside a Baseball?” and “What’s inside a Soccer Ball?” are in the same syntactic structure but the objectives are different. We define this as “text isomorphism.” Solely capitalizing on such titles will tend to draw the two videos close in representation space but unfortunately the videos describe two different objects, which are indicated by search queries. As a result, query information here is a rewarding signal to adjust visual-textual correlation and further improve video representation learning.

By delving into mutual calibration across query and text for weakly-supervised video representation learning, we present a novel Bi-Calibration Networks (BCN) architecture with the backbone of 3D ConvNet. Specifically, we employ the off-the-shelf BERT [12] model to extract the title features and cluster all the titles of the videos searched by an identical query into clusters. The centroid of each cluster is then taken as a text prototype and put into text vocabulary. Primary text supervision is measured as the cosine similarity between video title and all text prototypes in the vocabulary to regulate video-to-text (v2t) projection. Similarly, we take the “one-hot” vector in query vocabulary built directly on query words as primary query supervision to optimize video-to-query (v2q) projection. Next, v2t/v2q projection starts the text-to-query (t2q) or query-to-text (q2t) calibration, in which BCN aggregates/decomposes the predictions on cluster/query level in a bottom-up/top-down way to produce the t2q/q2t correction. The two corrections refine the primary query/text supervision to further optimize v2q/v2t projection. Moreover, we devise a selection scheme to balance the two corrections. The whole architecture is optimized by minimizing query and text classification loss.

The main contribution of this work is the proposal of exploring the cross correction between query and text to boost weakly-supervised video representation learning. This also leads to a better view of why query and text could complement to each other to validate visual-textual connections, and how to integrate the correction across the two into video representation learning framework. Two large-scale web video datasets are proposed for the weakly-supervised learning and extensive experiments on the datasets demonstrate the effectiveness of our BCN framework.

2 RELATED WORK

We briefly group the related works into three categories: supervised, unsupervised and weakly-supervised video representation learning with respect to the utilization of clean labels, no label or noisy labels for model training.

The early works [3], [13], [14], [15], [16], [17] of supervised video representation learning are extended from image representation by applying 2D CNN on video frames. For instance, Karpathy *et al.* [15] leverage spatio-temporal convolutions for representation learning on the stacked

frame-level features. To further capture the motion information, the well-known two-stream architecture [3] and its variants [13], [17], [18] are devised by executing 2D CNN on optical flow. Though the methods improve the representation learning by formulating motion pattern, performing 2D CNN on video frames still limits the capacity of modeling long-range temporal dynamics. To alleviate this problem, LSTM-RNN [16] captures the long-term dependencies in videos by utilizing a long-short term memory (LSTM) auto-encoder. To treat the video clip as a temporal evolution unit rather than a sequence of independent frames, 3D CNN structures [1], [19], [20], [21], [22] are proposed to boost video feature learning on large-scale supervised video datasets [1], [23]. More recently, the video transformer works [24], [25], [26] are extended from the vision transformers [27], [28] in image domain to learn video representation. Note that the backbone of our BCN is 3D ConvNet but the representation is learnt in a weakly-supervised manner.

Unsupervised video representation learning is one kind of technique to employ unlabeled videos for representation learning. The related works leverage various supervision from the video data itself to build pretext tasks for video representation learning, such as frame/clip order prediction [29], [30], [31], [32], motion estimation [33], [34], temporal cycle-consistency learning [35], [36], temporal coherence learning [37], [38], pixel-level displacement prediction [36], [39], frame reconstruction [40], [41], [42] and contrastive learning [43], [44]. To further improve the descriptive ability of video representation, weakly-supervised methods [10], [11], [45] focus on utilizing the weak supervision from web video data [46]. For example, Ghadiyaram *et al.* [10] explore the influences of different aspects in tags (e.g., label space) for weakly-supervised learning. Furthermore, to mine the knowledge from title of web videos, CPD [11] learns the video representation by making visual-textual pair close to each other. Nevertheless, most of the recent weakly-supervised methods are still facing the challenge of query polysemy or text isomorphism when directly exploiting query or text as supervision. Particularly, some previous works [47], [48], [49] handled query polysemy in web data collection, however our proposal alleviates it in the stage of representation learning.

In short, our approach belongs to weakly-supervised video representation learning techniques. Different from the aforementioned methods which solely rely on the supervision of query or text, our BCN in this paper contributes by studying not only mining supervisory signal from query and text simultaneously, but also how mutual calibration between query and text information could be leveraged to enhance weakly-supervised video representation learning.

3 BI-CALIBRATION NETWORKS

In this section, we introduce the Bi-Calibration Networks (BCN) that performs mutual calibration between query and text to facilitate weakly-supervised video representation learning. Specifically, BCN formulates the problem as the classification over query and text vocabulary. To better understand the spirit of our BCN design, we first introduce the problem formulation of the unique pretext task in BCN

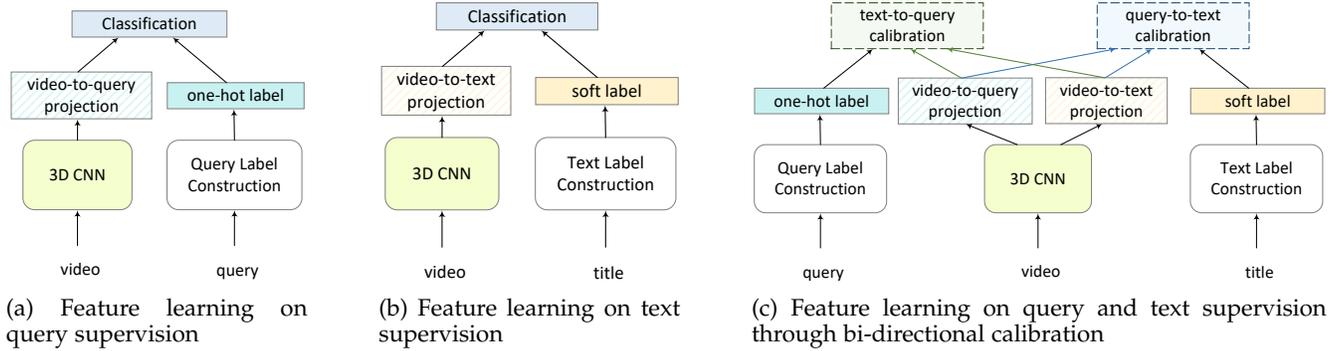


Fig. 2: The conceptual illustration of weakly-supervised video representation learning on (a) query supervision; (b) text supervision; (c) both of query and text supervision with bi-calibration in BCN.

for weakly-supervised video representation learning. After that, the detailed BCN architecture is further elaborated.

3.1 Problem Formulation

Suppose we have a web video collection, where each video is paired with the searched query and the video title description. The goal of the weakly-supervised video representation learning is to pre-train a visual encoder on the video data supervised by the query or text supervision. The pre-trained visual encoder can be further utilized to support several downstream tasks. In this paper, we employ the 3D CNN as the visual encoder for video pre-training.

For the video pre-training base on query, the query words are usually converted to “one-hot” labels. One natural way [10] to optimize the 3D CNN is to classify the video-to-query (v2q) projection of the 3D backbone based on these semantic labels, as shown in Figure 2(a). On the other hand, there are many directions to employ text to supervise video representation learning, such as the contrastive learning [11], triplet ranking [50] or linear feature regression [51]. Instead, we still formulate the weakly-supervised video feature learning on text as the classification pretext task. As illustrated by Figure 2(b), we convert the video titles into the soft labels and further employ them as the text supervision to optimize the 3D CNN through the classification on video-to-text (v2t) projection. To improve the quality of the primary query/text supervision and facilitate video representation learning, we design the bi-directional calibration learning paradigm to correct query and text supervision across each. The text-to-query (t2q) or query-to-text (q2t) calibration takes both of the v2q and v2t projections as the correction signals to refine the primary “one-hot” query label or “soft” text label, as illustrated by Figure 2(c).

Figure 3 further details an overview of the BCN framework and two coupled calibration modules, i.e., text-to-query (t2q) and query-to-text (q2t) calibration. Specifically, for each input video, BCN first utilizes a 3D CNN to extract video representation and feed it into two branches, i.e., video-to-query (v2q) and video-to-text (v2t) projections. Two kinds of probability distributions (i.e., the query distribution over query vocabulary and the text distribution over all text prototypes) are achieved to trigger each calibration. Note that all the titles of the videos searched by an identical query are initially grouped into multiple clusters, and each text prototype corresponds to the centroid of

each cluster. In this way, we naturally obtain two kinds of primary supervision for each video, i.e., the primary query supervision (the “one-hot” vector in query vocabulary) and the primary text supervision (the cosine similarity between video title and all text prototypes), that are used to optimize v2q and v2t projections, respectively. After that, BCN starts t2q/q2t calibration by aggregating/decomposing the text/query distribution into t2q/q2t correction in a bottom-up/top-down fashion, respectively. The learnt t2q/q2t correction is further integrated with the primary query/text supervision, yielding the refined query/text supervision to strengthen the regulation of each branch. During training, a selection scheme is utilized to balance the two calibrations.

3.2 Video-to-Query/Text Projection Branch

The ultimate target of BCN is to train the 3D CNN backbone for video representation learning through two pretext tasks of query and text classification. Therefore, based on the extracted video feature, we involve two parallel video-to-query (v2q) and video-to-text (v2t) projection branches to perform the two pretext tasks, which can be optimized with the corresponding query and text supervision.

Primary Query and Text Supervision. Given all the training web videos paired with the searched queries and surrounding text (i.e., titles), one natural way to represent each query or text as query/text supervision is to directly construct query/text vocabulary based on query/title words. Therefore, for each query, we adopt this natural way to take its “one-hot” vector $\mathbf{y}^q \in \mathbb{R}^K$ in query vocabulary (vocabulary size: K) as the primary query supervision. However, the robustness of this recipe is brittle when applying it to represent surrounding texts, since the inherently semantic distribution among different titles is unexploited. Inspired by the classical bag-of-words paradigm for feature representation [53], we leverage such paradigm to construct text vocabulary by delving into the diverse semantic structure of titles. Formally, we first extract the text feature \mathbf{f}^t of each video title by the off-the-shelf language model BERT [12]. Next, for the k -th query, we perform clustering on all the titles of the videos searched by this query through k -means algorithm, leading to m_k clusters. The cluster number m_k is automatically set by a statistic-based clustering estimation method (Gap Statistic [54]). The centroid \mathbf{B}_i of the i -th cluster is then defined as one text prototype, which is computed as the average of all text features of

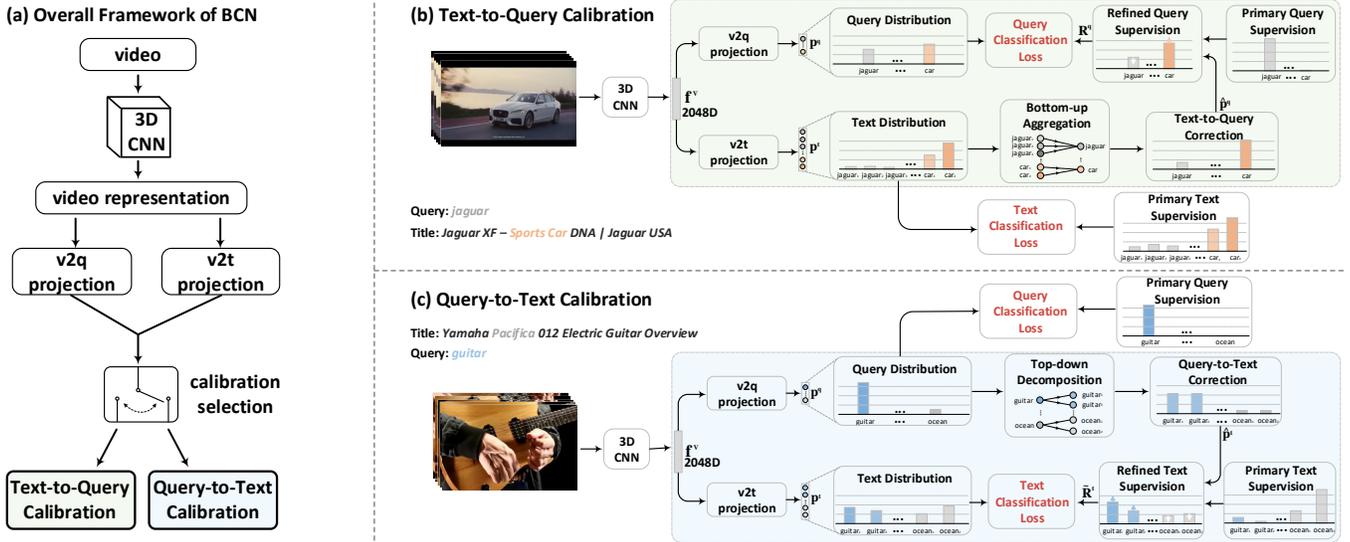


Fig. 3: An overview of our Bi-Calibration Networks (BCN). In general, BCN (a) first extracts the video representation of input video via 3D CNN, and then feed it into two branches, i.e., video-to-query (v2q) and video-to-text (v2t) projections. The output query distribution over query vocabulary and the text distribution over all text prototypes are utilized to trigger the text-to-query (b) and query-to-text (c) calibration modules, which are controlled by a selection scheme.

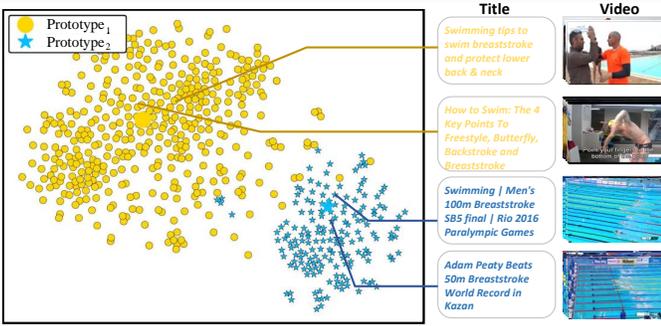


Fig. 4: Visualization of clustering on title features by t-SNE [52]. (query: swimming breast stroke)

video titles in that cluster. Accordingly, we construct the text vocabulary based on all the $M = \sum_{k=0}^{K-1} m_k$ text prototypes derived from K queries, which naturally characterize the underlying semantic structure of titles. Based on this text vocabulary, we interpret the primary text supervision $\mathbf{y}^t = [y_0^t, y_1^t, \dots, y_{M-1}^t] \in \mathbb{R}^M$ of each surrounding text \mathbf{f}^t as the cosine similarities between \mathbf{f}^t and all the M text prototypes $\{\mathbf{B}_i | i = 0, 1, \dots, M-1\}$:

$$\mathbf{y}^t = \text{softmax}(\text{COS}(\mathbf{f}^t, \mathbf{B}_i) | i = 0, 1, \dots, M-1), \quad (1)$$

where $\text{softmax}(\cdot)$ is the softmax with temperature parameter $1/M$ and $\text{COS}(\cdot)$ represents cosine similarity.

One case of the clustering for video titles searched by the query of “swimming breast stroke” is visualized in Figure 4. There are two meaningful clusters as shown in the figure. The first cluster has an emphasis on the daily breaststroke (e.g., tutorial) while most of the titles in the second cluster tend to describe swimeet in the championship. The differences in between are reflected in visual content. Thus, the text supervision based on the similarity between title feature and each prototype well captures the semantic relation.

Video-to-Query Projection. The v2q projection branch is especially designed to perform the pretext task of query classification. Concretely, depending on the input video feature \mathbf{f}^v extracted by 3D CNN, the v2q projection branch infers its query distribution $\mathbf{p}^q = [p_0^q, p_1^q, \dots, p_{K-1}^q]$ over query vocabulary through a fully-connected layer. Note that p_k^q is the estimated probability of assigning this video to the k -th query word. The v2q projection branch can be directly optimized with the primary query supervision \mathbf{y}^q , and we measure the query classification loss as softmax loss:

$$L^q = - \sum_{k=0}^{K-1} I_{k=y^q} \log(p_k^q), \quad (2)$$

where the indicator function $I_{k=y^q} = 1$ if the label value of k -th query word in \mathbf{y}^q is 1, otherwise $I_{k=y^q} = 0$.

Video-to-Text Projection. In analogy to the v2q projection branch, we design another v2t projection branch to conduct the pretext task of text classification. In particular, the v2t projection branch takes the video feature \mathbf{f}^v as the input, and learns to estimate the text distribution $\mathbf{p}^t = [p_0^t, p_1^t, \dots, p_{M-1}^t]$ over all text prototypes in text vocabulary. p_i^t denotes the probability of the video belonging to the i -th cluster. Accordingly, we optimize the v2t projection branch with the primary text supervision $\mathbf{y}^t = [y_0^t, y_1^t, \dots, y_{M-1}^t]$, and the text classification loss for this proxy task is calculated as the softmax cross-entropy loss:

$$L^t = - \sum_{i=0}^{M-1} y_i^t \log(p_i^t). \quad (3)$$

3.3 Text-to-Query Calibration

The most typical way to optimize video-to-query projection branch is to use the primary query supervision for query classification as in Eq.(2). However, such way oversimplifies the proxy task by assuming that all videos searched by an

identical query belong to one class, while ignoring the phenomenon of query polysemy (i.e., the coexistence of many possible meanings for a query). That will inevitably mislead video representation learning. To alleviate this problem, we devise a text-to-query (t2q) calibration module that further regulates the v2q projection with additional calibration from video surrounding texts (i.e., titles).

Text-to-Query Correction. The t2q calibration module first transforms the text distribution \mathbf{p}^t derived from the v2t projection branch into the t2q correction $\hat{\mathbf{p}}^q \in \mathbb{R}^K$ by aggregating all the probabilities of text prototypes belonging to the same query in a bottom-up way:

$$\hat{p}_k^q = \sum_{i \in I_k} p_i^t, \quad s.t. \quad k \in \{0, 1, \dots, K-1\}, \quad (4)$$

where \hat{p}_k^q denotes the k -th correction value in t2q correction $\hat{\mathbf{p}}^q$ and represents the aggregated probability with regard to the k -th query word in query vocabulary. I_k is the index set of the text prototypes for the k -th query.

Refined Query Supervision. Next, the t2q correction $\hat{\mathbf{p}}^q$ serves as the additional calibration from the text distribution over all text prototypes, aiming to refine the primary query supervision \mathbf{y}^q with more semantic meanings derived from v2t projection branch. Specifically, we first take the primary query supervision \mathbf{y}^q weighted by the estimated query distribution \mathbf{p}^q as the query confidence $\mathbf{y}^{q*} = \mathbf{y}^q \circ \mathbf{p}^q$, which reflects the confidence score of each query word. Note that the operation \circ denotes the element-wise multiplication. The refined query supervision \mathbf{R}^q is then estimated by integrating the t2q correction $\hat{\mathbf{p}}^q$ with the query confidence score \mathbf{y}^{q*} as:

$$\mathbf{R}^q = (\mathbf{y}^{q*} + \hat{\mathbf{p}}^q) / \|\mathbf{y}^{q*} + \hat{\mathbf{p}}^q\|_1, \quad (5)$$

where $\|\cdot\|_1$ represents \mathcal{L}_1 -norm. The underlying assumption behind Eq.(5) is that if the confidence score of the ground-truth query word is higher than the correction value of that query, the refined query supervision prefers to be closer to the primary query supervision. Otherwise, the refined query supervision tends to be heavily influenced with the t2q correction. Finally, the t2q calibration module leverages the refined query supervision \mathbf{R}^q to further optimize the v2q projection branch, and the query classification loss in Eq.(2) is thus reformulated as:

$$\hat{L}^q = - \sum_{k=0}^{K-1} R_k^q \log(p_k^q), \quad (6)$$

where R_k^q is the updated ground-truth of k -th query word in the refined query supervision \mathbf{R}^q .

3.4 Query-to-Text Calibration

Recall that the aforementioned optimization of the v2t projection branch (see Eq.(3)) solely hinges on the primary text supervision for text classification proxy task. Nevertheless, in the case of text isomorphism (i.e., titles share the same syntactic structure but refer to different semantics), the semantic discriminativeness of video representations learnt in this way may be easily overwhelmed. To address the issue, a query-to-text (q2t) calibration module is designed to guide the optimization of v2t projection with the additional high-level semantic supervision from query.

Query-to-Text Correction. In the q2t calibration module, we first calculate the q2t correction $\hat{\mathbf{p}}^t \in \mathbb{R}^M$ by evenly decomposing each element (e.g., the probability p_k^q of k -th query word) of query distribution to the correction values over all the text prototypes belonging to that query:

$$\hat{p}_i^t = \frac{1}{m_k} p_k^q, \quad s.t. \quad i \in I_k, \quad k \in \{0, 1, \dots, K-1\}, \quad (7)$$

where \hat{p}_i^t denotes the correction value of the i -th text prototype in q2t correction $\hat{\mathbf{p}}^t$ and m_k is the number of text prototypes belonging to the k -th query.

Refined Text Supervision. After that, we improve the primary text supervision with the q2t correction, leading to the refined text supervision to further regulate the v2t projection branch. In particular, the text confidence $\mathbf{y}^{t*} = \mathbf{y}^t \circ \mathbf{p}^t$ is first calculated as the primary text supervision \mathbf{y}^t weighted by text distribution \mathbf{p}^t . Each element in \mathbf{y}^{t*} denotes the confidence score of each text prototype. Similar to the formulation of refined query supervision, we measure the refined text supervision \mathbf{R}^t by fusing the q2t correction $\hat{\mathbf{p}}^t$ and the text confidence \mathbf{y}^{t*} :

$$\mathbf{R}^t = (\mathbf{y}^{t*} + \hat{\mathbf{p}}^t) / \|\mathbf{y}^{t*} + \hat{\mathbf{p}}^t\|_1. \quad (8)$$

Moreover, in order to make the refined text supervision evolve smoothly, we utilize a moving average update strategy with momentum to update the refined text supervision:

$$\tilde{\mathbf{R}}^t \leftarrow \alpha \tilde{\mathbf{R}}^t + (1 - \alpha) \mathbf{R}^t, \quad (9)$$

where α is a momentum coefficient, and \mathbf{R}^t or $\tilde{\mathbf{R}}^t$ denotes the new or running value of the refined text supervision, respectively. Based on the updated refined text supervision $\tilde{\mathbf{R}}^t$, the q2t calibration module further optimizes the v2t projection branch through the text classification loss:

$$\hat{L}^t = - \sum_{i=0}^{M-1} \tilde{R}_i^t \log(p_i^t), \quad (10)$$

where \tilde{R}_i^t is the new ground-truth of the i -th text prototype in the updated refined text supervision $\tilde{\mathbf{R}}^t$.

3.5 Network Optimization

During training, we adopt a two-stage strategy to optimize the whole architecture of our BCN. In the first stage, we optimize the BCN framework with the typical query and text classification losses (L^q in Eq.(2) and L^t in Eq.(3)) simultaneously, irrespective of any calibration modules. In the second stage, the BCN framework is further fine-tuned with two coupled calibration modules. Here, we design a selection scheme to balance the two calibration modules according to the difference between the query/text confidence ($\mathbf{y}^{q*}/\mathbf{y}^{t*}$) and t2q/q2t correction ($\hat{\mathbf{p}}^q/\hat{\mathbf{p}}^t$) in each module. Most specifically, if $\|\mathbf{y}^{q*} - \hat{\mathbf{p}}^q\|_2 > \varepsilon^q$ and $\|\mathbf{y}^{t*} - \hat{\mathbf{p}}^t\|_2 < \varepsilon^t$ (ε^q and ε^t are two thresholds), this case implies that the primary text supervision (\mathbf{y}^t) is more reliable than the primary query supervision (\mathbf{y}^q). Therefore we select the t2q calibration module, and the BCN framework is optimized with the modified query classification loss \hat{L}^q in Eq.(6) plus the typical text classification loss L^t in Eq.(3). In contrast, if $\|\mathbf{y}^{q*} - \hat{\mathbf{p}}^q\|_2 < \varepsilon^q$ and $\|\mathbf{y}^{t*} - \hat{\mathbf{p}}^t\|_2 > \varepsilon^t$, the primary query supervision is supposed to be more reliable than the primary text supervision, and the q2t calibration module is

TABLE 1: The statistics of YOVO-3M and YOVO-10M datasets.

| Dataset | Source | Supervision | # of Video | # of Clip | # of Query | Duration (hrs) |
|----------|----------|-----------------|------------|------------|------------|----------------|
| YOVO-3M | web data | query and title | 986,031 | 2,958,092 | 2,015 | 8216.9 |
| YOVO-10M | web data | query and title | 8,051,431 | 10,023,532 | 18,305 | 12142.1 |

Algorithm 1: Calibration Selection Scheme

Input : Model \hat{M} of the first training stage, web videos with query and title; two thresholds ε^q and ε^t ;

Output: Output model \hat{M} of BCN;

- 1 Initialize BCN model with \hat{M} , the iterative count $n = 0$;
- 2 **while** $n \leq N - 1$ **do**
- 3 Network forward to obtain probabilities \mathbf{p}^q and \mathbf{p}^t ;
- 4 Compute $\hat{\mathbf{p}}^q$ and $\hat{\mathbf{p}}^t$ according to Eq.(4) and Eq.(7);
- 5 Compute \mathbf{y}^{q*} , \mathbf{R}^q and \mathbf{y}^{t*} , $\hat{\mathbf{R}}^t$ by Eq.(5) and Eq.(8), Eq.(9);
- 6 **if** $\|\mathbf{y}^{q*} - \hat{\mathbf{p}}^q\|_2 > \varepsilon^q$ **and** $\|\mathbf{y}^{t*} - \hat{\mathbf{p}}^t\|_2 < \varepsilon^t$ **then**
- 7 Optimize \hat{M} by \hat{L}^q and L^t ; (t2q calibration)
- 8 **else if** $\|\mathbf{y}^{t*} - \hat{\mathbf{p}}^t\|_2 > \varepsilon^t$ **and** $\|\mathbf{y}^{q*} - \hat{\mathbf{p}}^q\|_2 < \varepsilon^q$ **then**
- 9 Optimize \hat{M} by \hat{L}^t and L^q ; (q2t calibration)
- 10 **else** Optimize \hat{M} by L^q and L^t ;
- 11 $n = n + 1$;
- 12 **end while**
- 13 **return** \hat{M}

selected for optimization. In that case, the objective of BCN framework consists of the typical query classification loss L^q in Eq.(2) and the modified text classification loss \hat{L}^t in Eq.(10). Otherwise, we optimize BCN with the typical query and text classification losses as in the first stage. The weight for each loss is set as 1.0 empirically.

Algorithm 1 details the processing of the selection scheme of the two calibrations in our BCN.

4 EXPERIMENTS

The experiments of weakly-supervised video representation learning are conducted on the newly-created YouTube video datasets, namely YOVO-3M and YOVO-10M, particularly paired with query and title. We then empirically verify the merit of BCN on three scene-related action recognition datasets: Kinetics-400 [1], UCF101 [55] and HMDB51 [56], and two interaction-related action recognition datasets: Something-Something V1 and V2 [57].

4.1 Datasets

YOVO-3M/10M Datasets. We collect the YOVO-3M and YOVO-10M datasets characterized by the unique properties including large-scale web video data with query and title information, as well as the comprehensive and diverse video content for weakly-supervised video representation learning. Figure 5 depicts the construction pipeline of YOVO-3M and YOVO-10M datasets, which consists of four main steps, i.e., query vocabulary collection, query deduplication, video collection on YouTube and clip deduplication. To crawl the web videos with comprehensive visual content, we first collect all the labels from Kinetics-400 [1], Kinetics-700 [58], ImageNet [59] and Moments [60] datasets as search queries. After query deduplication, the number of remaining queries is 2,015. We issue each query to YouTube and about 489

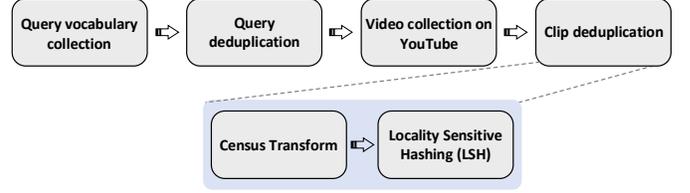


Fig. 5: The construction pipeline of YOVO-3M and YOVO-10M.

videos with the titles are downloaded successfully on average. We then uniformly sample three 10-second clips from each video to build YOVO-3M dataset. To further enlarge the volume of web videos, a 3K verb vocabulary and a 13K noun vocabulary extracted from Oxford English Dictionary are utilized as additional queries to search for another set of video clips with the titles. We combine these clips with YOVO-3M dataset to construct YOVO-10M dataset. Finally, we employ the standard clip deduplication approach [10] to remove video clips occurring anywhere in the downstream datasets from both YOVO-3M and YOVO-10M datasets. Specifically, we extract the global feature of each frame in a video clip through Census Transform [61] and then execute Locality Sensitive Hashing (LSH) [62] on the feature to obtain frame-level hash codes. We then compute the hamming distance between the frame-level hash codes of each pair of frames, in which one is from a video clip from YOVO-3M/10M datasets and the other is from a video in downstream datasets, e.g., Kinetics-400. We average all the distances of frame pairs across two videos as the clip-level distance in between. If the distance is lower than 2, we will remove the corresponding video clip from YOVO-3M/10M.

Table 1 summarizes the statistics of YOVO-3M and YOVO-10M datasets. In detail, YOVO-3M contains 2,015 queries and 2,958,092 video clips in total, and YOVO-10M consists of 18,305 queries and 10,023,532 video clips. The two scales of the proposed datasets would be applicable to different researchers w.r.t computational resources in the video representation learning community. Figure 6 further illustrates 32 video clips with the searched queries and titles from YOVO-3M and YOVO-10M datasets. The showing video cases demonstrate the diverse video content in different facets, e.g., objects, sports and daily activities, for weakly-supervised video representation learning.

Downstream Datasets. We evaluate our BCN on five downstream datasets. The Kinetics-400 [1] dataset consists of 300K 10-second clips from 400 action categories. There are 240K, 20K, 40K clips in training, validation and testing sets, respectively. The UCF101 [55] contains 13K videos from 101 action classes, and the HMDB51 [56] has 7K videos from 51 action categories. In UCF101 and HMDB51 datasets, there are three training/validation splits provided by the dataset organizers. Each split in UCF101 includes about 9.5K training and 3.7K validation videos, and a HMDB51

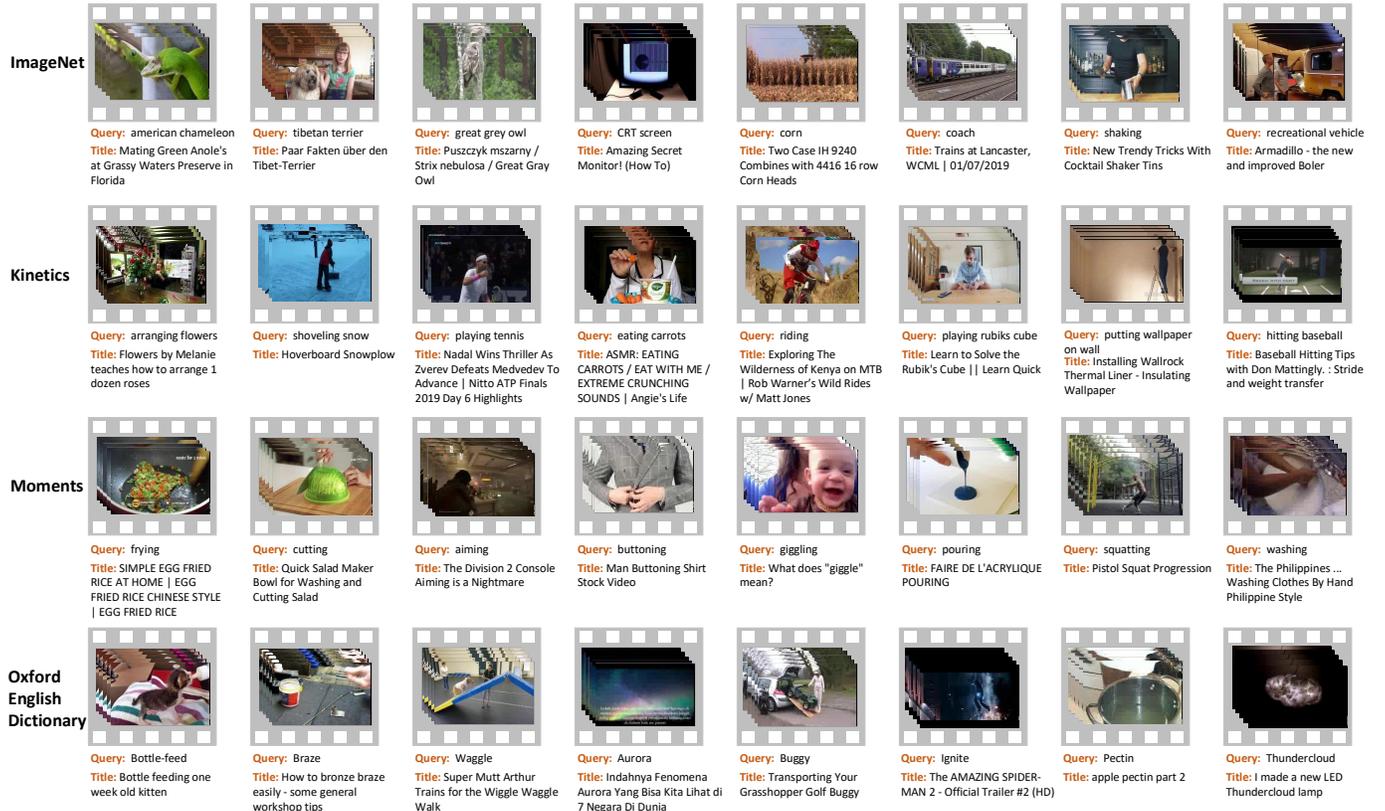


Fig. 6: Examples of 32 video clips with searched queries and titles from YOVO-3M and YOVO-10M datasets. The video clips in each row are searched by the queries in the category list of different datasets, i.e., ImageNet, Kinetics and Moments, or in the vocabulary list of the Oxford English Dictionary.

split contains 3.5K training and 1.5K validation videos. The videos in the above three datasets are mainly for scene-related action recognition. In Something-Something V1 and V2 [57] datasets, there are about 108K and 221K videos from 174 action categories, respectively. The training/validation/testing set consists of 86K/11.5K/11K and 169K/25K/27K videos, which are mostly for interaction-related recognition.

4.2 Experimental Settings

Implementations. We employ BERT [12] text encoder to extract the 1,024-D features for each token in the title and average all the token features as text representation. The number M of total clusters is 6,819/65,766 on YOVO-3M/10M. In weakly-supervised training, we utilize the architecture of LGD-3D [63] based on ResNet-50 [64] as the network backbone but all the parameters are trained from scratch. The dimension of the input video clips is set as $32 \times 112 \times 112$, which is randomly cropped from the original web video. Each clip is randomly flipped along horizontal direction for data augmentation. We choose the two threshold ε^q and ε^t as 0.5 and 0.7 by cross validation. The momentum coefficient α is fixed to 0.9. We implement BCN on Caffe [65] platform. In all the pre-training stages, the networks are trained by utilizing stochastic gradient descent (SGD) with 0.9 momentum. The initial learning rate is set to 0.08 and 0.008 in the first and second training stage, and decreased by 10% after every 200K iterations. The mini-batch size is set as 256 and the weight decay is 0.0001.

Evaluation Metrics. We adopt two evaluation protocols in the downstream datasets, i.e., linear model and network fine-tuning. In the former protocol, we uniformly sample 10 or 3 video clips from each video in Kinetics-400/UCF101/HMDB51 or Something-Something V1/V2 datasets, and take the 2,048-way outputs from pool5 layer of the network backbone as the features of each clip. We average all the features of clips in one video as video representation, and a linear SVM is learnt on the training set and evaluated on validation set. In detail, the cost parameter c in SVM is set as 8.0. Both top-1 and top-5 accuracy are reported as evaluation metrics. In the latter one, we initialize the network backbone with the weakly-supervised training output model of BCN, and fine-tune/evaluate the network on the training/validation set of each dataset.

4.3 Evaluation on Primary Text Supervision

We first examine the effectiveness of primary text supervision for video representation learning, regardless of mutual calibration design. We compare the following four methods: (1) The regression method (RG) optimizes video representation through minimizing Smooth L1 loss [66] between video representation and text feature. (2) Triplet Ranking algorithm (TR) learns video representation to make positive video-text pair more similar than negative pair. (3) A variant of our primary text supervision (TS-) also performs clustering on video titles and builds text vocabulary on the centroids of clusters. Each title is naturally assigned to one cluster and represented as a binary index vector in text

TABLE 2: Top-1 and Top-5 accuracy on Kinetics-400, UCF101 and HMDB51 under linear protocol. (Training on YOYO-3M).

| Approach | Kinetics-400 | | UCF101 | | HMDB51 | |
|----------|--------------|-------|--------|-------|--------|-------|
| | Top-1 | Top-5 | Top-1 | Top-5 | Top-1 | Top-5 |
| RG | 64.5 | 84.9 | 85.1 | 94.3 | 58.3 | 81.2 |
| TR | 67.2 | 87.2 | 87.2 | 95.7 | 61.8 | 83.9 |
| TS- | 71.5 | 89.1 | 91.0 | 98.9 | 62.2 | 85.6 |
| TS | 72.5 | 89.6 | 91.8 | 99.0 | 64.5 | 89.7 |

TABLE 3: Performance contribution of each design in BCN. The model is learnt on YOYO-3M and evaluated by linear protocol. The ‘‘Query’’ and ‘‘Text’’ denote the corresponding supervision.

| Method | Query | Text | T→Q | Q→T | K400 | U101 | HD51 |
|------------------|-------|------|-----|-----|-------------|-------------|-------------|
| QS | ✓ | | | | 71.1 | 90.7 | 62.1 |
| TS | | ✓ | | | 72.5 | 91.8 | 64.5 |
| QS+TS | ✓ | ✓ | | | 73.1 | 92.3 | 65.9 |
| BCN _Q | ✓ | ✓ | ✓ | | 73.8 | 93.3 | 66.9 |
| BCN _T | ✓ | ✓ | | ✓ | 73.6 | 92.8 | 66.3 |
| BCN | ✓ | ✓ | ✓ | ✓ | 74.2 | 93.5 | 67.6 |

vocabulary. We exploit single-label classification to regulate video representation with text supervision in TS-. (4) TS is our proposed primary text supervision in BCN, which computes cosine similarity between title and all clusters. TS can be regarded as the soft mode of TS- and accordingly tunes video feature with multi-label classification loss.

Table 2 summarizes performance comparisons of video representation learnt with different ways of primary text supervision under linear protocol on three downstream datasets. Overall, the results across different datasets consistently indicate that TS leads to a performance boost against other methods. In particular, the top-1 accuracy of our TS achieves 72.5%, 91.8% and 64.5% on Kinetics, UCF101 and HMDB51, respectively, making the absolute improvement over RG/TR by 8.0%/5.3%, 6.7%/4.6% and 6.2%/2.7%. Such results demonstrate the advantage of exploring the structure among all the titles of videos via clustering. Though both TS- and TS utilize text clustering to improve text supervision, they are different in the way that TS- represents each title as an index vector (1 for its own cluster, otherwise 0), and TS is by computing cosine similarity between the title and all clusters. As indicated by the results, delving into the correlation between each title and all title clusters in TS leads to better performances.

Discussion with Contrastive Learning. Inspired from the self-supervised learning [67], the video representation learning from text is recently formulated as the contrastive learning [11] between the visual and textual features. The features of the positive video-text pairs are pulled close while the features of negative pairs are pushed away. Following such setting, we experimented with contrastive learning on video and text via InfoNCE loss [68] on the same backbone. The top-1 accuracy achieves 72.2% on Kinetics-400 under linear evaluation protocol. The performance is lower than 72.5% of the TS in Table 2. Instead of learning the pair-wise correlation through contrastive learning, our TS mines the group-wise relationship via clustering. The results basically confirm that our TS is a good alternative for the video representation learning based on text information.

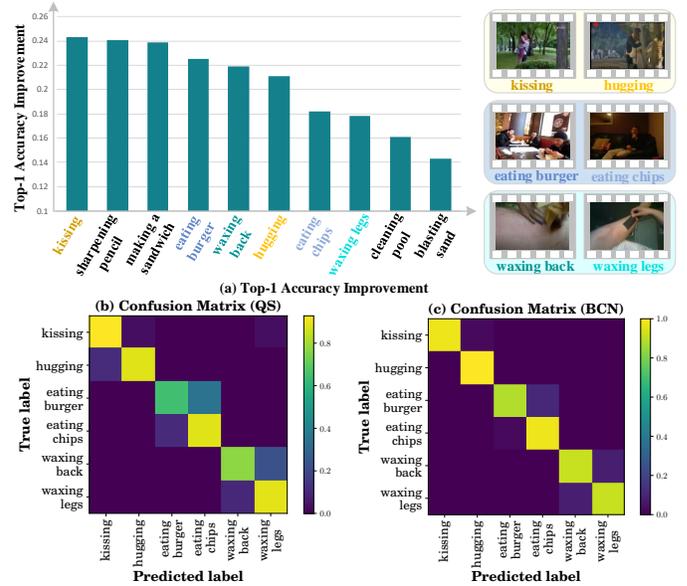


Fig. 7: (a) 10 classes with the largest performance gain from QS to BCN under linear protocol on Kinetics-400 and the visualization of confusion matrix across the selected six categories of video representation learnt by (b) QS and (c) BCN.

4.4 Evaluation on Bi-Calibration Networks

Next, we study how each design in Bi-Calibration Networks influences video representation learning. QS and TS solely exploit primary query supervision and text supervision in our framework to guide video representation learning, respectively, through single-label and multi-label classification. QS+TS performs the joint learning on the primary supervision of query and text. BCN_Q/BCN_T leverages the idea of text-to-query (T→Q) or query-to-text (Q→T) calibration to estimate t2q/q2t corrections to adjust primary query/text supervision and further boost video feature learning. BCN is our Bi-Calibration framework.

Table 3 details the top-1 accuracy on Kinetics-400 (K400), UCF101 (U101) and HMDB51 (HD51) datasets under linear model protocol by considering one more factor in our BCN. Compared to QS/TS, QS+TS boosts up the accuracy from 71.1%/72.5%, 90.7%/91.8% and 62.1%/64.5%, to 73.1%, 92.3% and 65.9%, respectively, on three datasets. The results basically indicate that primary query supervision and primary text supervision are complementary to refine visual-textual connections and enhance video representation learning. Text-to-query/query-to-text calibration further amends primary query/text supervision and the performance gain of each is 0.7%/0.5%, 1.0%/0.5% and 1.0%/0.4% on three datasets against QS+TS. By mutual calibration between query and text, BCN finally reaches the top-1 accuracy of 74.2%, 93.5% and 67.6%.

To verify the impact of mutual calibration via our BCN design across different categories, we further list the categories with most benefit. Figure 7(a) shows ten categories in Kinetics-400 which achieve the largest performance gain from QS to BCN. An interesting observation is that three pairs, i.e., kissing-hugging, eating burger-eating chips, waxing back-waxing legs, among the ten categories, are indeed fine-grained and it is challenging to distinguish one from the other in each pair. Figure 7(b) and (c) also visualizes

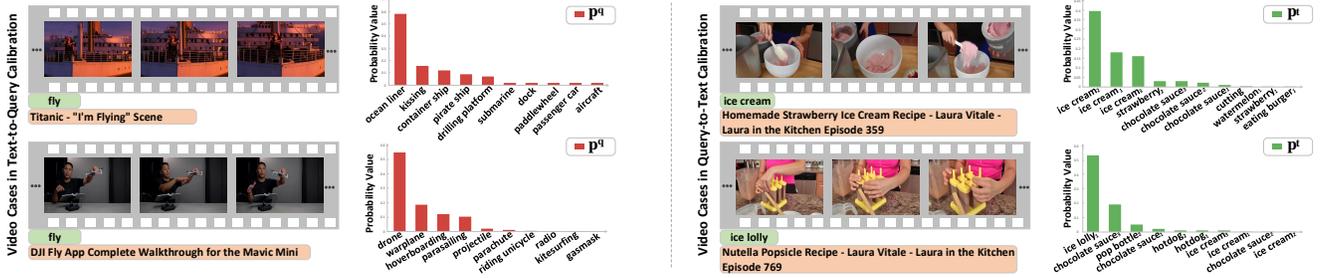


Fig. 8: Video examples with the searched queries, video titles and the probability distribution on query/text (p^q/p^t).

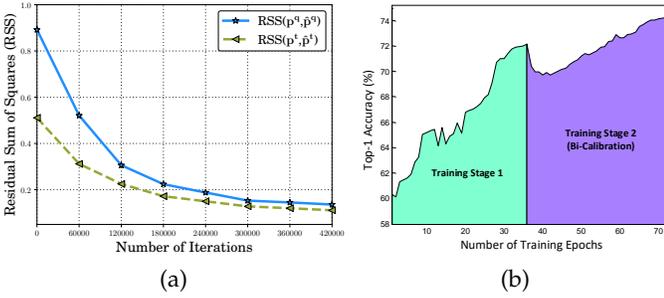


Fig. 9: (a) Residual Sum of Square (RSS) curves between p^q/p^t and \hat{p}^q/\hat{p}^t w.r.t the number of iterations. (b) Top-1 accuracy under linear protocol of two training stages on Kinetics-400. (The number of training epoch in each stage is 36 and decided through cross validation.)

confusion matrix across the six categories on video representation learnt by QS and BCN, respectively. It is clear that BCN endows video representation with more discriminative power especially on the fine-grained categories.

4.5 More Analysis on Network Optimization

The selection scheme controls the switch across the calibration of text-to-query or query-to-text directions. We compute the Residual Sum of Squares (RSS) between the probabilities (p^q/p^t) and t2q/q2t corrections (\hat{p}^q/\hat{p}^t), and Figure 9(a) depicts RSS curve with respect to the number of iterations. As expected, the RSS on both (p^q, \hat{p}^q) and (p^t, \hat{p}^t) is gradually decreased when training more iterations. The curve of RSS on (p^t, \hat{p}^t) changes more smoothly due to the momentum update strategy. Since the corrections are adopted as the supervision to optimize the probabilities, the gradients are not back-propagated to them. Thus, the results give the clue that probabilities are enforced to be close to the corrections and BCN makes necessary corrections through text-to-query/query-to-text calibration, validating the impact of selection scheme. Figure 9(b) shows the curve of top-1 accuracy on Kinetics-400 under linear protocol in two optimization stages. Despite having some drops in accuracy at the beginning of bi-calibration training stage, the top-1 accuracy eventually improves and reaches a higher value, which again demonstrates the effectiveness of calibration between query and text through information refinement.

Figure 8 showcases two groups of videos with the searched queries, video titles and the probabilities on query/text. The first group of videos are searched by an identical query of “fly,” but the video content corresponds to different meanings. One is about the famous scene on

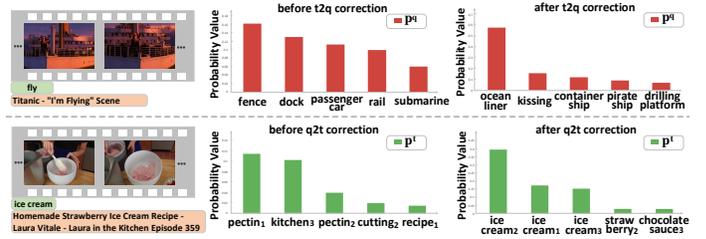


Fig. 10: Probability distribution on query/text (p^q/p^t) of the two cases in Fig 8 before and after the t2q/q2t correction.

the ocean liner in movie “Titanic” and the other describes “DJI” drone. Such query polysemy may misdirect video representation learning. With the text-to-query corrections in BCN, the query probabilities of the two videos are well predicted. The video about “Titanic” is highly relevant to “ocean liner” and “kissing,” and the video of “DJI” has high probability in response to “drone” and “warplane.” In contrast, the second group of videos share similar syntactic structure of title, but are related to different queries of “ice cream” and “ice lolly.” In this case, solely capitalizing on title information may affect video representation learning as well. Through query-to-text correction, the text probabilities predicted by BCN nicely lead to different emphasis of “ice cream” or “ice lolly” and “chocolate sauce.” Figure 10 further illustrates the probability distribution on query/text of two video cases in Figure 8 before and after the t2q/q2t correction. For the “Titanic” video, the top-2 query prediction before t2q correction are “fence” and “dock”. After the t2q calibration, we achieve more accurate predictions (“ocean liner” and “kissing”). Meanwhile, the q2t correction changes the original prediction of “pectin” to the preciser one of “ice cream” in the second video case.

4.6 Comparisons with State-of-the-Art Methods

We compare BCN with several state-of-the-art techniques on five datasets: Kinetics-400 (K400), UCF101 (U101) and HMDB51 (HD51) for scene-related action recognition, and Something-Something V1 (SS-V1) and V2 (SS-V2) for interaction-related action recognition.

Table 4 lists the top-1 and top-5 accuracy of different approaches on K400. For fair comparison, all the methods exploit only RGB modality of frame for model training and “*” denotes that the models are pre-trained on our YOYO-3M/10M with official source codes. Under linear model protocol, BCN achieves comparable performances with models built on fully-supervised ImageNet pre-training and self-

TABLE 4: Performance comparisons on Kinetics-400.

| Approach | Pre-training | Backbone | Top-1 | Top-5 |
|---|--------------|-----------|-------------|-------------|
| Supervised Pre-training | | | | |
| R(2+1)D RGB [22] | ImageNet | custom | 74.3 | 91.4 |
| I3D RGB [1] | ImageNet | Inception | 72.1 | 90.3 |
| S3D RGB [69] | ImageNet | Inception | 74.7 | 93.4 |
| NL I3D RGB [70] | ImageNet | ResNet-50 | 74.9 | 91.6 |
| TSM RGB [71] | ImageNet | ResNet-50 | 74.1 | 91.2 |
| LGD RGB [63] | ImageNet | ResNet-50 | 74.8 | 92.0 |
| SlowFast RGB [72] | ImageNet | ResNet-50 | 75.6 | 92.1 |
| SmallBig RGB [73] | ImageNet | ResNet-50 | 76.3 | 92.5 |
| TDN RGB [74] | ImageNet | ResNet-50 | 77.5 | 93.2 |
| Self-supervised Pre-training (linear protocol) | | | | |
| VTHCL RGB [75] | K400 | ResNet-50 | 37.8 | - |
| CVRL RGB [76] | K400 | ResNet-50 | 66.1 | - |
| ρ BYOL RGB [43] | K400 | ResNet-50 | 71.5 | - |
| Weakly-supervised Pre-training | | | | |
| Linear model protocol on video representation | | | | |
| CPD RGB [11] | K400-title | ResNet-50 | 63.8 | - |
| BCN RGB | YOVO-3M | ResNet-50 | 74.2 | 90.9 |
| BCN RGB | YOVO-10M | ResNet-50 | 74.9 | 91.6 |
| Network fine-tuning | | | | |
| CPD* RGB [11] | YOVO-3M | ResNet-50 | 73.9 | 90.2 |
| CPD* RGB [11] | YOVO-10M | ResNet-50 | 75.0 | 91.3 |
| BCN RGB | YOVO-3M | ResNet-50 | 78.5 | 94.3 |
| BCN RGB | YOVO-10M | ResNet-50 | 79.1 | 94.7 |

supervised Kinetics pre-training. When fine-tuning the pre-trained model by BCN on K400, BCN exhibits better performances against other baselines. In particular, BCN learnt on YOVO-3M obtains 78.5% top-1 accuracy, which outperforms CPD [11] by 4.6% pre-trained on the same data. Different from CPD which solely capitalizes on title information, BCN exploits mutual calibration between query and title to improve weakly-supervised video representation learning. Such result demonstrates the advantage of our bi-calibration design. Compared to SmallBig and TDN, BCN leads the top-1 accuracy by 2.2% and 1.0%, respectively. Executing weakly-supervised learning on YOVO-10M further improves the accuracy from 78.5% to 79.1%. The performance trends on U101 and HD51 are similar with those on K400 as shown in Table 5. The results again verify the impact of BCN for weakly-supervised representation learning. Table 6 summarizes the performances on SS-V1 and SS-V2. Similarly, BCN under network fine-tuning protocol surpasses the best competitor SmallBig and ACTION-Net by 1.9% and 1.8% on SS-V1 and SS-V2, respectively.

5 CONCLUSIONS AND DISCUSSIONS

We have presented Bi-Calibration Networks (BCN), which explores the correlations between web videos and the searched queries or video titles for improving weakly-supervised video representation learning. Particularly, we study the problem from the viewpoint of refining the visual-semantic connections through mutual calibration between query and title information. To materialize our idea, we first achieve the primary query and text supervision on query vocabulary of query words and text vocabulary of text prototypes, which are utilized to optimize video-to-query (v2q) and video-to-text (v2t) projections for classification. Next,

TABLE 5: Top-1 accuracy on UCF101 and HMDB51.

| Approach | Pre-training | Backbone | U101 | HD51 |
|---|---------------|------------|-------------|-------------|
| Supervised Pre-training | | | | |
| R(2+1)D RGB [22] | K400 | custom | 96.8 | 74.5 |
| I3D RGB [1] | Img+K400 | Inception | 95.4 | 74.5 |
| S3D RGB [69] | Img+K400 | Inception | 96.8 | 75.9 |
| TSM RGB [71] | K400 | ResNet-50 | 95.9 | 73.5 |
| LGD RGB [63] | Img+K600 | ResNet-50 | 96.0 | 74.7 |
| Self-supervised Pre-training (fine-tuning) | | | | |
| XDC RGB [77] | K400 | R(2+1)D-18 | 84.2 | 47.1 |
| SpeedNet RGB [78] | K400 | S3D-G | 81.1 | 48.8 |
| CoCLR RGB [79] | K400 | S3D-G | 87.9 | 54.6 |
| CVRL RGB [76] | K400 | ResNet-50 | 92.2 | 66.7 |
| ρ BYOL RGB [43] | K400 | ResNet-50 | 95.5 | 73.6 |
| Weakly-supervised Pre-training | | | | |
| Linear model protocol on video representation | | | | |
| BCN RGB | YOVO-3M | ResNet-50 | 93.5 | 67.6 |
| BCN RGB | YOVO-10M | ResNet-50 | 94.9 | 69.2 |
| Network fine-tuning | | | | |
| MIL-NCE RGB [45] | HowTo100M | ResNet-50 | 91.3 | 61.0 |
| CPD RGB [11] | Instagram300k | ResNet-50 | 92.8 | 63.8 |
| CPD* RGB [11] | YOVO-3M | ResNet-50 | 94.2 | 70.2 |
| CPD* RGB [11] | YOVO-10M | ResNet-50 | 95.4 | 73.5 |
| BCN RGB | YOVO-3M | ResNet-50 | 97.1 | 74.9 |
| BCN RGB | YOVO-10M | ResNet-50 | 97.5 | 76.5 |

TABLE 6: Top-1 accuracy on Something-Something V1/V2.

| Approach | Pre-training | Backbone | SS-V1 | SS-V2 |
|---|---------------|-----------|-------------|-------------|
| Supervised Pre-training | | | | |
| I3D RGB [80] | ImageNet+K400 | ResNet-50 | 41.6 | - |
| NL I3D RGB [80] | ImageNet+K400 | ResNet-50 | 44.4 | - |
| NL I3D + gcn RGB [80] | ImageNet+K400 | ResNet-50 | 46.1 | - |
| CPNet RGB [81] | ImageNet | ResNet-34 | - | 57.7 |
| TSM RGB [71] | ImageNet | ResNet-50 | 45.6 | 59.1 |
| SmallBig RGB [73] | ImageNet | ResNet-50 | 48.3 | 61.6 |
| ACTION-Net RGB [82] | ImageNet | ResNet-50 | - | 62.5 |
| Self-supervised Pre-training (fine-tuning) | | | | |
| BYOL RGB [43] | K400 | ResNet-50 | - | 55.8 |
| MoCo RGB [43] | K400 | ResNet-50 | - | 54.4 |
| Weakly-supervised Pre-training | | | | |
| Linear model protocol on video representation | | | | |
| BCN RGB | YOVO-3M | ResNet-50 | 42.7 | 53.2 |
| BCN RGB | YOVO-10M | ResNet-50 | 43.2 | 55.7 |
| Network fine-tuning | | | | |
| CPD* RGB [11] | YOVO-3M | ResNet-50 | 45.2 | 60.1 |
| CPD* RGB [11] | YOVO-10M | ResNet-50 | 47.1 | 61.7 |
| BCN RGB | YOVO-3M | ResNet-50 | 48.6 | 62.6 |
| BCN RGB | YOVO-10M | ResNet-50 | 50.2 | 64.3 |

the v2t/v2q projection triggers the text-to-query or query-to-text calibration, that aims to adjust primary query/text supervision to further optimize v2q/v2t projection. Extensive experiments conducted on newly-created web video datasets, i.e., YOVO-3M and YOVO-10M, validate our BCN. More remarkably, weakly-supervised pre-training BCN on YOVO-10M is superior to several techniques with fully-supervised ImageNet or Kinetics pre-training.

REFERENCES

- [1] J. Carreira and A. Zisserman, "Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset," in *CVPR*, 2017. **1, 2, 6, 10**
- [2] Y. Li, B. Ji, X. Shi, J. Zhang, B. Kang, and L. Wang, "TEA: Temporal Excitation and Aggregation for Action Recognition," in *CVPR*, 2020. **1**
- [3] K. Simonyan and A. Zisserman, "Two-stream Convolutional Networks for Action Recognition in Videos," in *NIPS*, 2014. **1, 2**
- [4] C. Yang, Y. Xu, J. Shi, B. Dai, and B. Zhou, "Temporal Pyramid Network for Action Recognition," in *CVPR*, 2020. **1**

- [5] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu, "Video Swin Transformer," *arXiv preprint arXiv:2106.13230*, 2021. **1**
- [6] S. Ji, W. Xu, M. Yang, and K. Yu, "3D Convolutional Neural Networks for Human Action Recognition," *IEEE Trans. on PAMI*, vol. 35, no. 1, pp. 221–231, 2012. **1**
- [7] G. Varol, I. Laptev, and C. Schmid, "Long-Term Temporal Convolutions for Action Recognition," *IEEE Trans. on PAMI*, vol. 40, no. 6, pp. 1510–1517, 2017. **1**
- [8] X. Wang, L. Zhu, Y. Wu, and Y. Yang, "Symbiotic Attention for Egocentric Action Recognition with Object-centric Alignment," *IEEE Trans. on PAMI*, vol. 1, no. 1, pp. 1–13, 2020. **1**
- [9] A. Gaidon, Z. Harchaoui, and C. Schmid, "Temporal Localization of Actions with Actoms," *IEEE Trans. on PAMI*, vol. 35, no. 11, pp. 2782–2795, 2013. **1**
- [10] D. Ghadiyaram, M. Feiszli, D. Tran, X. Yan, H. Wang, and D. Mahajan, "Large-scale weakly-supervised pre-training for video action recognition," in *CVPR*, 2019. **1, 2, 3, 6**
- [11] T. Li and L. Wang, "Learning Spatiotemporal Features via Video and Text Pair Discrimination," *arXiv preprint arXiv:2001.05691*, 2020. **1, 2, 3, 8, 10**
- [12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *ACL*, 2018. **2, 3, 7**
- [13] A. Diba, V. Sharma, and L. V. Gool, "Deep Temporal Linear Encoding Networks," in *CVPR*, 2017. **2**
- [14] C. Feichtenhofer, A. Pinz, and A. Zisserman, "Convolutional Two-Stream Network Fusion for Video Action Recognition," in *CVPR*, 2016. **2**
- [15] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, "Large-scale Video Classification with Convolutional Neural Networks," in *CVPR*, 2014. **2**
- [16] J. Y.-H. Ng, M. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici, "Beyond Short Snippets: Deep Networks for Video Classification," in *CVPR*, 2015. **2**
- [17] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. V. Gool, "Temporal Segment Networks: Towards Good Practices for Deep Action Recognition," in *ECCV*, 2016. **2**
- [18] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Gool, "Temporal Segment Networks for Action Recognition in Videos," *IEEE Trans. on PAMI*, vol. 41, no. 11, pp. 2740–2755, 2018. **2**
- [19] K. Hara, H. Kataoka, and Y. Satoh, "Can Spatiotemporal 3D CNNs Retrace the History of 2D CNNs and ImageNet?" in *CVPR*, 2018. **2**
- [20] Z. Qiu, T. Yao, and T. Mei, "Learning Spatio-Temporal Representation with Pseudo-3D Residual Networks," in *ICCV*, 2017. **2**
- [21] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning Spatiotemporal Features with 3D Convolutional Networks," in *ICCV*, 2015. **2**
- [22] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri, "A Closer Look at Spatiotemporal Convolutions for Action Recognition," in *CVPR*, 2018. **2, 10**
- [23] B. Ghanem, J. C. Niebles, C. Snoek, F. C. Heilbron, H. Alwassel, V. Escorcia, R. Krishna, S. Buch, and C. D. Dao, "The ActivityNet Large-Scale Activity Recognition Challenge 2018 Summary," *arXiv preprint arXiv:1808.03766*, 2018. **2**
- [24] G. Bertasius, H. Wang, and L. Torresani, "Is Space-Time Attention All You Need for Video Understanding?" in *ICML*, 2021. **2**
- [25] H. Fan, B. Xiong, K. Mangalam, Y. Li, Z. Yan, J. Malik, and C. Feichtenhofer, "Multiscale Vision Transformers," *arXiv preprint arXiv:2104.11227*, 2021. **2**
- [26] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lucic, and C. Schmid, "ViViT: A Video Vision Transformer," in *ICCV*, 2021. **2**
- [27] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," in *ICLR*, 2021. **2**
- [28] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows," in *ICCV*, 2021. **2**
- [29] B. Fernando, H. Bilen, E. Gavves, and S. Gould, "Self-Supervised Video Representation Learning With Odd-One-Out Networks," in *CVPR*, 2017. **2**
- [30] I. Misra, C. L. Zitnick, and M. Hebert, "Shuffle and Learn: Unsupervised Learning using Temporal Order Verification," in *ECCV*, 2016. **2**
- [31] D. Wei, J. Lim, A. Zisserman, and W. T. Freeman, "Learning and Using the Arrow of Time," in *CVPR*, 2018. **2**
- [32] D. Xu, J. Xiao, Z. Zhao, J. Shao, D. Xie, and Y. Zhuang, "Self-supervised Spatiotemporal Learning via Video Clip Order Prediction," in *CVPR*, 2019. **2**
- [33] P. Agrawal, J. Carreira, and J. Malik, "Learning to See by Moving," in *ICCV*, 2015. **2**
- [34] D. Pathak, R. Girshick, P. Dollar, T. Darrell, and B. Hariharan, "Learning Features by Watching Objects Move," in *CVPR*, 2017. **2**
- [35] D. Dwibedi, Y. Aytar, J. Tompson, P. Sermanet, and A. Zisserman, "Temporal Cycle-Consistency Learning," in *CVPR*, 2019. **2**
- [36] X. Wang, A. Jabri, and A. A. Efros, "Learning Correspondence from the Cycle-Consistency of Time," in *CVPR*, 2019. **2**
- [37] H. Mobahi, R. Collobert, and J. Weston, "Deep Learning from Temporal Coherence in Video," in *ICML*, 2009. **2**
- [38] X. Wang and A. Gupta, "Unsupervised Learning of Visual Representations using Videos," in *ICCV*, 2015. **2**
- [39] Z. Liu, R. A. Yeh, X. Tang, Y. Liu, and A. Agarwala, "Video Frame Synthesis using Deep Voxel Flow," in *ICCV*, 2017. **2**
- [40] C. Finn, I. Goodfellow, and S. Levine, "Unsupervised Learning for Physical Interaction through Video Prediction," in *NIPS*, 2016. **2**
- [41] Z. Luo, B. Peng, D.-A. Huang, A. Alahi, and L. Fei-Fei, "Unsupervised Learning of Long-Term Motion Dynamics for Videos," in *CVPR*, 2017. **2**
- [42] N. Srivastava, E. Mansimov, and R. Salakhutdinov, "Unsupervised Learning of Video Representations using LSTMs," in *ICML*, 2015. **2**
- [43] C. Feichtenhofer, H. Fan, B. Xiong, R. Girshick, and K. He, "A Large-Scale Study on Unsupervised Spatiotemporal Representation Learning," in *CVPR*, 2021. **2, 10**
- [44] J. Li, P. Zhou, C. Xiong, and S. Hoi, "Prototypical Contrastive Learning of Unsupervised Representations," in *ICLR*, 2021. **2**
- [45] A. Miech, J.-B. Alayrac, L. Smaira, I. Laptev, J. Sivic, and A. Zisserman, "End-to-End Learning of Visual Representations from Uncurated Instructional Videos," in *CVPR*, 2020. **2, 10**
- [46] A. Miech, D. Zhukov, J.-B. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic, "HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips," in *ICCV*, 2019. **2**
- [47] T. L. Berg and D. A. Forsyth, "Animals on Web," in *CVPR*, 2006. **2**
- [48] K. Saenko and T. Darrell, "Unsupervised Learning of Visual Sense Models for Polysemous Words," in *NIPS*, 2008. **2**
- [49] F. Schroff, A. Criminisi, and A. Zisserman, "Harvesting Image Databases from the Web," in *ICCV*, 2007. **2**
- [50] J. C. Stroud, D. A. Ross, C. Sun, J. Deng, R. Sukthankar, and C. Schmid, "Learning Video Representations from Textual Web Supervision," *arXiv preprint arXiv:2007.14937*, 2020. **3**
- [51] Z. Qiu, T. Yao, C.-W. Ngo, X.-P. Zhang, D. Wu, and T. Mei, "Boosting Video Representation Learning with Multi-Faceted Integration," in *CVPR*, 2021. **3**
- [52] L. van der Maaten and G. Hinton, "Visualizing Data using t-SNE," *JMLR*, 2008. **4**
- [53] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories," in *CVPR*, 2006. **3**
- [54] R. Tibshirani, G. Walthers, and T. Hastie, "Estimating the number of cluster in a data set via the gap statistic," *Journal of the Royal Statistical Society: Series B*, pp. 411–423, 2001. **3**
- [55] K. Soomro, A. R. Zamir, and M. Shah, "UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild," *CRCV-TR-12-01*, 2012. **6**
- [56] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre, "HMDB: A Large Video Database for Human Motion Recognition," in *ICCV*, 2011. **6**
- [57] R. Goyal, S. E. Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Freund, P. Yianilos, M. Mueller-Freitag, F. Hoppe, C. Thureau, I. Bax, and R. Memisevic, "The "something something" video database for learning and evaluating visual common sense," in *ICCV*, 2017. **6, 7**
- [58] J. Carreira, E. Noland, C. Hillier, and A. Zisserman, "A Short Note on the Kinetics-700 Human Action Dataset," *arXiv preprint arXiv:1907.06987*, 2019. **6**
- [59] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *CVPR*, 2009. **6**

- [60] M. Monfort, A. Andonian, B. Zhou, K. Ramakrishnan, S. A. Bargal, T. Yan, L. Brown, Q. Fan, D. Gutfreund, C. Vondrick, and A. Oliva, "Moments in Time Dataset: One Million Videos for Event Understanding," *IEEE Trans. on PAMI*, vol. 42, no. 2, pp. 502–508, 2019. [6](#)
- [61] R. Zabih and J. Woodfill, "Non-parametric Local Transforms for Computing Visual Correspondence," in *ECCV*, 1994. [6](#)
- [62] A. Gionis, P. Indyk, and R. Motwani, "Similarity Search in High Dimensions via Hashing," in *VLDB*, 1999. [6](#)
- [63] Z. Qiu, T. Yao, C.-W. Ngo, X. Tian, and T. Mei, "Learning Spatio-Temporal Representation with Local and Global Diffusion," in *CVPR*, 2019. [7](#), [10](#)
- [64] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *CVPR*, 2016. [7](#)
- [65] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. B. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional Architecture for Fast Feature Embedding," in *ACM MM*, 2014. [7](#)
- [66] R. Girshick, "Fast R-CNN," in *ICCV*, 2015. [7](#)
- [67] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum Contrast for Unsupervised Visual Representation Learning," in *CVPR*, 2020. [8](#)
- [68] A. van den Oord, Y. Li, and O. Vinyals, "Representation Learning with Contrastive Predictive Coding," in *NIPS*, 2018. [8](#)
- [69] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy, "Rethinking Spatiotemporal Feature Learning: Speed-Accuracy Trade-offs in Video Classification," in *ECCV*, 2018. [10](#)
- [70] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local Neural Networks," in *CVPR*, 2018. [10](#)
- [71] J. Lin, C. Gan, and S. Han, "TSM: Temporal Shift Module for Efficient Video Understanding," in *ICCV*, 2019. [10](#)
- [72] C. Feichtenhofer, H. Fan, J. Malik, and K. He, "SlowFast Networks for Video Recognition," in *ICCV*, 2019. [10](#)
- [73] X. Li, Y. Wang, Z. Zhou, and Y. Qiao, "SmallBigNet: Integrating Core and Contextual Views for Video Classification," in *CVPR*, 2020. [10](#)
- [74] L. Wang, Z. Tong, B. Ji, and G. Wu, "TDN: Temporal Difference Networks for Efficient Action Recognition," in *CVPR*, 2021. [10](#)
- [75] C. Yang, Y. Xu, B. Dai, and B. Zhou, "Video Representation Learning with Visual Tempo Consistency," *arXiv preprint arXiv:2006.15489*, 2020. [10](#)
- [76] R. Qian, T. Meng, B. Gong, M.-H. Yang, H. Wang, S. Belongie, and Y. Cui, "Spatiotemporal Contrastive Video Representation Learning," in *CVPR*, 2021. [10](#)
- [77] H. Alwassel, D. Mahajan, B. Korbar, L. Torresani, B. Ghanem, and D. Tran, "Self-Supervised Learning by Cross-Modal Audio-Video Clustering," in *NeurIPS*, 2020. [10](#)
- [78] S. Benaim, A. Ephrat, O. Lang, I. Mosseri, W. T. Freeman, M. Rubinstein, M. Irani, and T. Dekel, "SpeedNet: Learning the Speediness in Videos," in *CVPR*, 2020. [10](#)
- [79] T. Han, W. Xie, and A. Zisserman, "Self-supervised Co-training for Video Representation Learning," in *NeurIPS*, 2020. [10](#)
- [80] X. Wang and A. Gupta, "Videos as Space-Time Region Graphs," in *ECCV*, 2018. [10](#)
- [81] X. Liu, J.-Y. Lee, and H. Jin, "Learning Video Representations from Correspondence Proposals," in *CVPR*, 2019. [10](#)
- [82] Z. Wang, Q. She, and A. Smolic, "ACTION-Net: Multipath Excitation for Action Recognition," in *CVPR*, 2021. [10](#)



Fuchen Long is currently a Researcher in Vision and Multimedia Lab at JD Explore Academy, Beijing, China. He has participated several temporal action proposal and detection competitions such as Activity detection in Extended Videos (ActEV-PC) in ActivityNet Challenge 2019, ActivityNet Temporal Action Detection Challenge 2018 and ActivityNet Temporal Action Proposal Challenge 2017. His research interests include temporal action proposal and localization, multimedia retrieval and video understanding. He received Ph.D. degree from the University of Science and Technology of China (USTC) in 2021.



Ting Yao is currently a Principal Researcher in Vision and Multimedia Lab at JD Explore Academy, Beijing, China. His research interests include video understanding, vision and language, and deep learning. Prior to joining JD.com, he was a Researcher with Microsoft Research Asia, Beijing, China. Ting is the principal designer of several top-performing multimedia analytic systems in international benchmark competitions such as ActivityNet Large Scale Activity Recognition Challenge 2019-2016, Visual Domain Adaptation Challenge 2019-2017, and COCO Image Captioning Challenge. He is the leader organizer of MSR Video to Language Challenge in ACM Multimedia 2017 & 2016, and built MSR-VTT, a large-scale video to text dataset that is widely used worldwide. His works have led to many awards, including ACM SIGMM Outstanding Ph.D. Thesis Award 2015, ACM SIGMM Rising Star Award 2019, and IEEE TCMC Rising Star Award 2019. He is also an Associate Editor of *IEEE Trans. on Multimedia*.



Zhaofan Qiu is currently a Researcher in Vision and Multimedia Lab at JD Explore Academy, Beijing, China. His research interests include large-scale video classification, semantic segmentation, and multimedia understanding. He has participated several large-scale video analysis competitions such as ActivityNet Large Scale Activity Recognition Challenge, and THUMOS Action Recognition Challenge. He was awarded the MSRA Fellowship in 2017. He received Ph.D. degree in 2020 from the University of Science and Technology of China (USTC), Hefei, China.



Xinmei Tian (M'13) is an Associate Professor in the CAS Key Laboratory of Technology in Geospatial Information Processing and Application System, University of Science and Technology of China. She received the B.E. degree and Ph.D. degree from the University of Science and Technology of China in 2005 and 2010, respectively. Her current research interests include multimedia information retrieval and machine learning. She received the Excellent Doctoral Dissertation of Chinese Academy of Sciences award in 2012 and the Nomination of National Excellent Doctoral Dissertation award in 2013.



Jiebo Luo (S93, M96, SM99, F09) joined the Department of Computer Science at the University of Rochester in 2011, after a prolific career of over 15 years with Kodak Research. He has authored over 400 technical papers and holds over 90 U.S. patents. His research interests include computer vision, machine learning, data mining, social media, and biomedical informatics. He has served as the Program Chair of the ACM Multimedia 2010, IEEE CVPR 2012, ACM ICMR 2016, and IEEE ICIP 2017, and on the Editorial Boards of the *IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE*, *IEEE TRANSACTIONS ON MULTIMEDIA*, *IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY*, *IEEE TRANSACTIONS ON BIG DATA*, *Pattern Recognition*, *Machine Vision and Applications*, and *ACM Transactions on Intelligent Systems and Technology*. He is also a Fellow of ACM, AAAI, SPIE and IAPR.



Tao Mei (M07-SM11-F19) is a Vice President with JD.COM and the Deputy Managing Director of JD Explore Academy, where he also serves as the Director of Computer Vision and Multimedia Lab. Prior to joining JD.COM in 2018, he was a Senior Research Manager with Microsoft Research Asia in Beijing, China. He has authored or co-authored over 200 publications (with 12 best paper awards) in journals and conferences, 10 book chapters, and edited five books. He holds over 25 US and international patents. He

is or has been an Editorial Board Member of IEEE Trans. on Image Processing, IEEE Trans. on Circuits and Systems for Video Technology, IEEE Trans. on Multimedia, ACM Trans. on Multimedia Computing, Communications, and Applications, Pattern Recognition, etc. He is the General Co-chair of IEEE ICME 2019, the Program Co-chair of ACM Multimedia 2018, IEEE ICME 2015 and IEEE MMSP 2015.

Tao received B.E. and Ph.D. degrees from the University of Science and Technology of China, Hefei, China, in 2001 and 2006, respectively. He is a Fellow of IEEE (2019), a Fellow of IAPR (2016), a Distinguished Scientist of ACM (2016), and a Distinguished Industry Speaker of IEEE Signal Processing Society (2017).