# Fine-Grained Instance-Level Sketch-Based Video Retrieval

Peng Xu, Kun Liu, Tao Xiang, Timothy M. Hospedales, Zhanyu Ma, Jun Guo, and Yi-Zhe Song

*Abstract*—Existing sketch-analysis work studies sketches depicting static objects or scenes. In this work, we propose a novel cross-modal retrieval problem of fine-grained instance-level sketch-based video retrieval (FG-SBVR), where a sketch sequence is used as a query to retrieve a specific target video instance. Compared with sketch-based still image retrieval, and coarse-grained category-level video retrieval, this is more challenging as both visual appearance and motion need to be simultaneously matched at a fine-grained level. We contribute the first FG-SBVR dataset with rich annotations. We then introduce a novel multi-stream multi-modality deep network to perform FG-SBVR under both strong and weakly supervised settings. The key component of the network is a relation module, designed to prevent model overfitting given scarce training data. We show that this model significantly outperforms a number of existing state-of-the-art models designed for video analysis.

*Index Terms*—fine-grained video retrieval, sketch-based video retrieval, sketch dataset, cross-modal matching, triplet ranking, meta-learning inspired techniques.



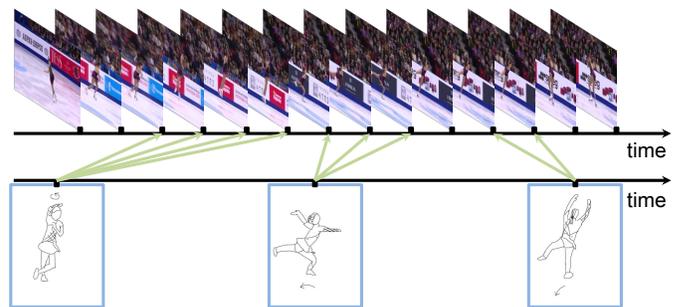Fig. 1. Sketch examples of the existing SBVR dataset [14].



Fig. 2. Illustration of fine-grained instance-level sketch-based video retrieval. A sketch-sequence (bottom) is connected to the video frames they summarise (top). Best viewed in color.

## I. Introduction

IT is said that one sketch speaks for a hundred words. Sketch provides a convenient abstraction to bridge concepts and pixels, via capturing both salient detail and topology. Sketch is now convenient and widely captured given the modern prevalence of touch-screen devices. This has led to a flourishing [1] of sketch-related research in recent years including sketch-based image retrieval [2], [3], [4], sketch-generation [5], segmentation [6], hashing [7], abstraction [8], scene understanding [9], [10], and self-supervised representation learning [11]. However, all of these studies only work with sketches depicting static objects or scenes, and analysis of sketching of motion is much under-studied.

Humans recall and describe events from episodic memory [12] with selective effects [13]. Visual recollections mainly contains the appearance and actions of key objects (*e.g.*, movements, spinning, rising). Combined with free-hand drawing of arrows or lines, sketch can simultaneously describe the appearance and motion of objects corresponding to such typical human recollections. Motivated by this potential, sketch-based video retrieval (SBVR) was first proposed in [14], and

followed up in several subsequent studies [15], [16], [17], [18]. However, these early methods are relatively *coarse-grained*, with sketched objects providing almost symbolic category indicators, rather than fine-grained details where sketch really shines as an alternative to conventional tagging approaches [19], [2]. Moreover the associated datasets are not large enough to train contemporary deep methods, and are not instance-level, in the sense of there being a set or videos rather than a single target video for each query sketch. Some examples of the TSF dataset introduced in [14] are shown in Figure 1, in which objects are *iconic*, without any fine-grained appearance information. This fails to exploit the full expressiveness of sketch and undermines the practical motivation for SBVR, since conventional symbolic tags ('person') could be a more convenient query modality there.

In this paper we provide the first study of genuinely fine-grained instance-level sketch-based video retrieval (FG-SBVR). This task is extremely challenging since it not only needs to solve all the difficulties common to static-sketch cross-modal retrieval (*i.e.*, matching abstract and sparse line-drawings to dense pixel renderings of perspective projections), but also requires understanding of motion depiction in sketch, and registering sketches to specific time windows within a temporally extended video. To support research in this area, we introduce the first FG-SBVR dataset containing $1,448$ sketches corresponding to $528$ figure skating video

Peng Xu is with School of Computer Science and Engineering, Nanyang Technological University, Singapore. E-mail: peng.xu@ntu.edu.sg Homepage: http://www.pengxu.net/ GitHub: https://github.com/PengBoXiangShang

Kun Liu, Zhanyu Ma, and Jun Guo are with Beijing University of Posts and Telecommunications, China. E-mail: {liu_kun, mazhanyu, guojun}@bupt.edu.cn

Tao Xiang, Yi-Zhe Song are with Centre for Vision, Speech and Signal Processing (CVSSP), University of Surrey, United Kingdom. E-mail: {t.xiang, y.song}@surrey.ac.uk

Timothy M. Hospedales is with University of Edinburgh, United Kingdom. E-mail: t.hospedales@ed.ac.uk

clips. Going beyond previous studies: (i) All sketches depict subtle appearance and pose details of skaters including body posture, hand gesture, clothing, and hair style. (ii) Skater movements are summarised by fine-grained motion vectors indicating skater glissades, spins and jumping. (iii) To represent temporally extended video rather than instantaneous actions, video clips are described by multi-page sketches. One example is illustrated in Figure 2. This multi-page "skater + motion vector" sketch format means that a successful FG-SBVR model must solve three challenging matching/alignment issues: (i) fine-grained visual appearance matching between drawn skaters and video image frames, (ii) fine-grained motion matching between sketched motion indicators and the video motion, (iii) alignment of pages within the sketch query to sub-sequences within the video.

In order to solve this cross-modal matching problem, we introduce the first deep model for FG-SBVR. It is a multi-stream multi-modality deep neural network. Specifically, the network is designed to align the video and sketch-sequence modalities in a joint embedding space so that they can be compared by a specific distance. Taking into consideration the unique "skater + motion vector" sketch format, each modality is modelled by a sub-network composed of a "appearance" stream and a "motion" stream. Within each stream, cross-modal matching is modelled by triplet ranking. This network design is against the recent trend in video analysis dominated by 3D convolutional networks [20], [21], [22], which do not explicitly separate dynamic and static video content. We found that with the large modality gap between video and sketch and the scarce training data, decomposing the dynamic and static aspects of both modalities explicitly becomes crucial.

To further address the training data scarcity problem, inspired by existing meta learning based few-shot learning work [23], [24], we introduce a relation module into our FG-SBVR. One of the most effective ways of improving generalization by meta learning is training a non-linear comparison – or relation – module. The relation module improves the learned representation by modelling the non-linear relationship between sketch-clip pairs, and benefits from learning from more negative pairs compared to triplet loss alone. Our FG-SBVR can be thus trained effectively even with sparse data. Furthermore, we explore both the strongly supervised setting (using ground-truth sketch page-frame alignment annotation during training), and the weakly-supervised learning (no within-video sketch-frame correspondence) setting based on multi-instance learning (MIL).

Our main contributions can be summarised as: (i) We propose the novel and challenging problem of fine-grained sketch-based video retrieval. (ii) We contribute the first FG-SBVR dataset with extensive ground truth annotation[1]. (iii) We develop a novel multi-stream multi-modality deep network to solve FG-SBVR by explicitly decomposing appearance and motion. A relation module is also introduced in the network to prevent overfitting. (iv) We explore learning this framework with both strong- and weak supervision using multi-instance learning. Extensive experiments are conducted to show that

the proposed model outperforms a number of state-of-the-art video analysis baselines.

## II. RELATED WORK

**Video Retrieval**  Many video retrieval techniques are query-by-example (QBE) [25], in which users provide (visual, textual, audio, *etc*) examples of the content they seek. According to query modality, video retrieval spans [26], [27]: (i) image-to-video (I2V) retrieval [27], (ii) text-to-video (T2V) retrieval [28], [29], and (iii) video-to-video (V2V) retrieval, *e.g.*, via hashing [30], [31], [32], [33], [34]. V2V is a unimodal task, T2V is a cross-modal task, and I2V is somewhere between. However, all these query modalities have some drawbacks: (i) Image query only provides appearance for one moment without dynamic clues. (ii) Text generally needs a lot of words or sentences (for example, the appearance, movement and relative position of objects) if the video is to be described in detail, rather than just used as a list of keyword tags. (iii) Video query matches the modality to the data to retrieve, but it may be difficult to obtain a representative video that matches the desired video to be retrieved.

**Sketch-Based Video Retrieval**  The pioneering work [14] on SBVR proposed a probabilistic model for content-based video retrieval driven by free-hand sketch queries depicting both objects and their movement (via dynamic cues, *i.e.*, lines and arrows). This work led to a series of subsequent research including tracking visual key-points in videos to form short trajectories which can be clustered to form tokens summarising video content. These can then be matched to a color and motion description of a query sketch by a Viterbi-like process [15]. This method was improved in [16] as a hybrid "semantic sketch" based video retrieval system by fusing the semantics of text with the expressiveness of sketch. Similarly, a Markov Random Field (MRF) optimisation based SBVR approach is proposed in [17], which combines shape, motion, colour and semantics within a single SBVR framework. Then, an index-based hybrid SBVR system is proposed in [18]. However, the mentioned SBVR approaches have several drawbacks: (i) Their retrievals are relatively coarse-grained (see Figure 1). This undermines the unique practical advantage of sketch: to convey cues that are hard to describe with simple symbolic tags [19], [2]. Existing video tagging systems already allow text/tag-based video search for such coarse concepts. (ii) Besides studying simple sketches, existing SBVR methods address retrieving relatively instantaneous actions (see Figure 1). They do not address the challenge of retrieving actions with temporally extended structure.

**SBVR Datasets**  There are only a few coarse-grained SBVR datasets [14] to date. As mentioned earlier, in addition to the lack of *instance-level* sketch-video pairing, all sketches of these datasets are overly iconic object contours with motion lines and arrows. To facilitate FG-SBVR research, we therefore introduce the first fine-grained SBVR dataset. The unique features of our dataset are: (i) It contains fine-grained instance-level SBVR data with one ground-truth video match to each sketch. (ii) It supports evaluation of more complex temporal logic in SBVR, by allowing more than one page of

---

[1]Our dataset and code will be made public.

action to describe a whole temporally extended video. (iii) It contains fine-grained visual detail enabling pose, clothing, and hairstyle to be used as matching cues.

**Video Analysis Models** Video analysis (*e.g.*, action recognition) methods mainly include RNN plus CNN [35], [36], two-stream networks [37], [38], and 3D convolutional networks (*e.g.*, C3D [39], P3D [40], I3D [41], T-C3D [22], 3D ResNet18 [21], ARTNets [42], Non-Local Neural Network [20]), STC-Net [43], S3D [44], and MFNet [45]). As mentioned earlier, we argue that the two-stream architecture is particularly suited for our multi-steam multi-modality alignment task because the dynamic and static streams of the video naturally correspond to the motion vector and static skater parts of sketches, respectively. This is validated in our experiments (see Section V-B). Our model is related to fine-grained sketch-based image retrieval (FG-SBIR) models [2], [19], but deals with a more challenging multi-stream matching problem. Importantly, none of the existing video analysis or FG-SBIR models exploit a relation module to address the overfitting problem due to data scarcity.

## III. Fine-Grained Instance-Level SBVR Dataset

We contribute the first fine-grained instance-level sketch-based video retrieval dataset. It contains 528 HD figure skating clips and 1,448 corresponding sketches. Our sketches contain both fine-grained appearance information, local dynamics, and longer-time frame dynamics via sequences of sketch 'pages'. Some examples are shown in Figure 3. Compared to sketches in prior datasets (Figure 1), ours have significantly more fine-grained details. We next describe the data collection and annotation process, and provide some quantitative comparisons with prior datasets.

### A. Data Collection

**Videos** We download diverse professional figure skating competition videos (*e.g.*, US National, European and World Championships) from YouTube. From these, we selected 49 female figure skating videos (duration: 6 to 56 minutes). For each video, both 720P and 1080P files are stored at 30 FPS including audio channels with English narratives. The audio channel can support future research in speech or text modalities (*e.g.*, extracting the keywords from narratives as 'attribute vectors' describing the video). We recruited 5 skating fans to select representative clips from the original long videos. We cut out 528 clips, with a total duration of 3,546 seconds. The average length is 6.7 seconds, with minimum 1 and maximum 29 seconds. Detailed duration statistics are shown in Figure 4(a).

**Sketches** The second step is to sketch the collected videos. We recruited 17 skating fans who are amateur sketchers to sketch the clips. As can be seen in Figure 3, due to lack of prior art training, these sketches are representative of the drawing abilities of the general population. The volunteers have a warm-up exercise, and then for each video clip, the volunteer can watch it several times and sketch what he/she has seen on a tablet, using their fingers or tablet stylus. Following recent practice [46], [19], our sketches are saved

TABLE I
COMPARISON WITH PREVIOUS SKETCH DATASETS. "-" MEANS NO STROKE INFORMATION.

| | Strokes | | | | Resolution ($W \times H$) | Sketch Amount |
|---|---|---|---|---|---|---|
| | min | max | mean | std | | |
| TU-Berlin [46] | 1 | 318 | 17.55 | 16.83 | $800 \times 800$ | 20k |
| Sketchy [19] | 1 | 434 | 17.91 | 16.06 | $640 \times 480$ | 75k |
| QMUL Shoe | - | - | - | - | $256 \times 256$ | 419 |
| QMUL chair | - | - | - | - | $256 \times 256$ | 297 |
| QMUL handbag | - | - | - | - | $256 \times 256$ | 568 |
| QMUL Chair-V2 | 1 | 138 | 12.79 | 9.84 | $800 \times 800$ | 1,275 |
| Ours | 26 | 345 | 102.40 | 43.47 | $768 \times 1024$ | 1,448 |

in Scalable Vector Graphic (SVG) format that stores spatio-temporal stroke information.

Each sketch contains two parts: the skater depicted at certain posture representing a key moment of the routine, and a motion vector summarising the movements of the skater centred around that key moment. The motion vector is abstract and subjective, so some instructions are necessary to avoid some completely random interpretations. In particular, the volunteers were told that: (i) If a video clip just shows a static scene (*e.g.*, skater keeps a static posture), draw the skater without any motion vectors. (ii) For jumping, the vector should be drawn above the skater's head. For gliding and spinning, the vector should be drawn below the skater's feet. (iii) Only one motion vector can be drawn within one sketch. (iv) One principle [47] of information processing theory states that short-term memory is organised as chunks of meaningful units. Thus, the volunteers are free to decide how-many pages of sketch to use to summarise the video clip.

Among the 1,448 sketches, there are 1,384 motion sketches (with motion vectors, containing different kinds of spinning, jumping, *etc*) and 64 static sketches (without motion vectors). The sketch sequences range from one to nine pages. Interestingly, 520 ($\approx 98\%$) of our video clips use less than seven sketches, which is consistent with human short term memory capacity being limited to seven chunks [47]. On average 2.7 sketches are used to describe each video clip, and detailed statistics are shown in Figure 4(b). In Figure 4(d), we present a scatter plot to show the relation between video duration and corresponding sketch sequence, in which the radius of each point is proportional to the number of video clips. We can see that people draw more sketches as video duration increases.

**Comparison with Other Datasets** Compared with existing sketch datasets, our dataset contains more details as reflected by the greater number of strokes (see Figure 4(c)). This comparison is made quantitatively in Table I. It shows that our dataset contains similar number of sketches as previous single category sketch based cross-modal retrieval datasets (*i.e.*, all except TU-Berlin [46] for sketch recognition and Sketchy [19] with 125 categories and 600 sketches per category). Drawing sketches with a reference image/video is very tedious making collecting large-scale datasets extremely difficult. Designing models that can be learned effectively with scarce data is thus a common challenge in sketch-based retrieval tasks.

Fig. 3. Examples of our figure skating FG-SBVR dataset. For each sketch page, its corresponding video frames (6 frames are selected) are shown.
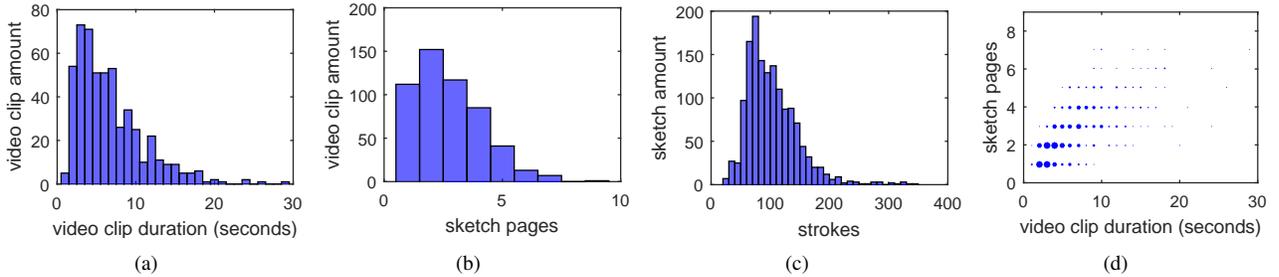


Fig. 4. Statistical analysis of our FG-SBVR dataset.

*B. Data Annotation*

**Motion Stroke Annotation** While drawing each sketch, the volunteers annotate which strokes form the motion vector. This simplifies separation of the "skater" and "motion" components of each sketch for our proposed multi-stream multi-modality model (see the sketch input part in Figure 5).

**Sketch-to-Frames Annotation** To provide strong supervision for correspondence between multi-page sketches and videos, the volunteers next annotate which video frames correspond to each of their sketches. This correspondence annotation is illustrated in Figure 2. We will explore the importance of using this information later (see Section V).

## IV. METHODOLOGY

**Problem Setting** We assume that the training dataset $D$ consists of $N$ paired sketch sequences and video clips: $D = \{(S_i, V_i)\}_{i=1}^N$. Each sketch sequence $S_i$ is composed of $M_i$ sketch 'pages', and each of these has an appearance and motion component: $S_i = \{(s_j^{ap}, s_j^{mo})\}_{j=1}^{M_i}$. Similarly, each video clip $V_i$ is composed of $O_i$ frame chunks: $V_i = \{(v_j^{ap}, v_j^{mo})\}_{j=1}^{O_i}$. Given $D$, we aim to learn a deep sketch and video mult-modal joint embedding space, where the similarity of a sketch query and video pair can be simply computed as a distance for retrieval.
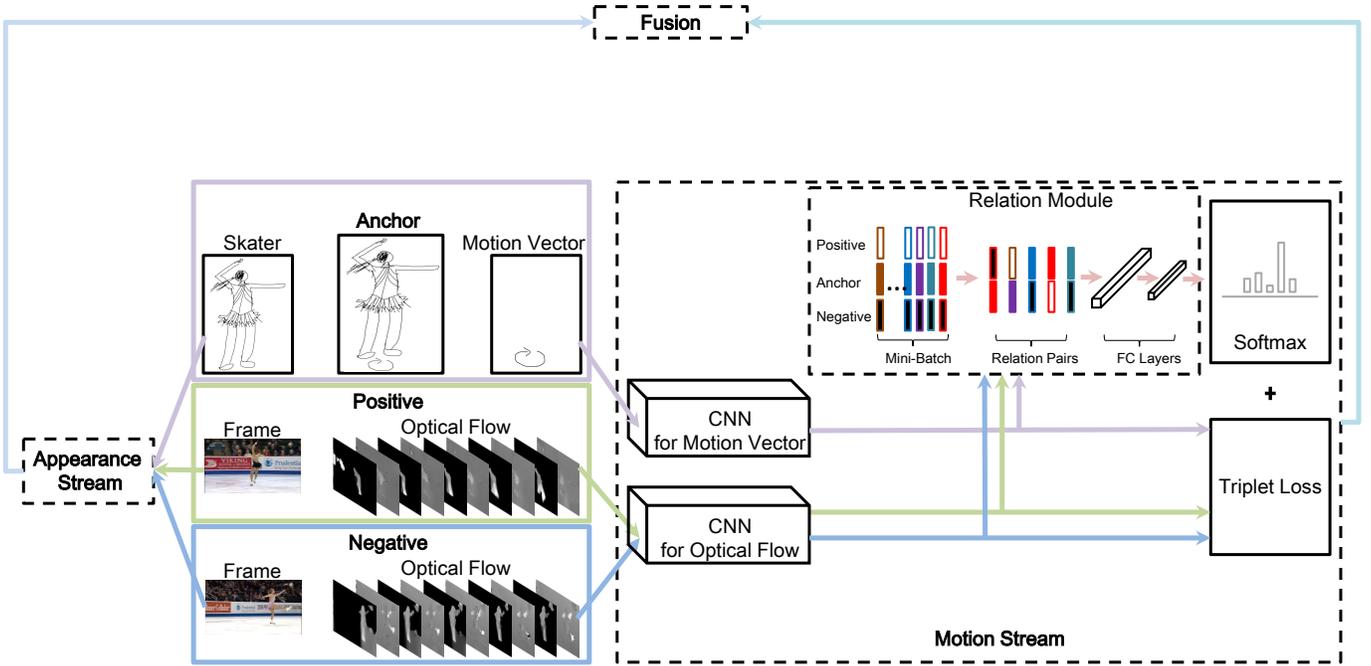
Fig. 5. Architecture of the proposed framework for FG-SBVR. Best viewed in color.
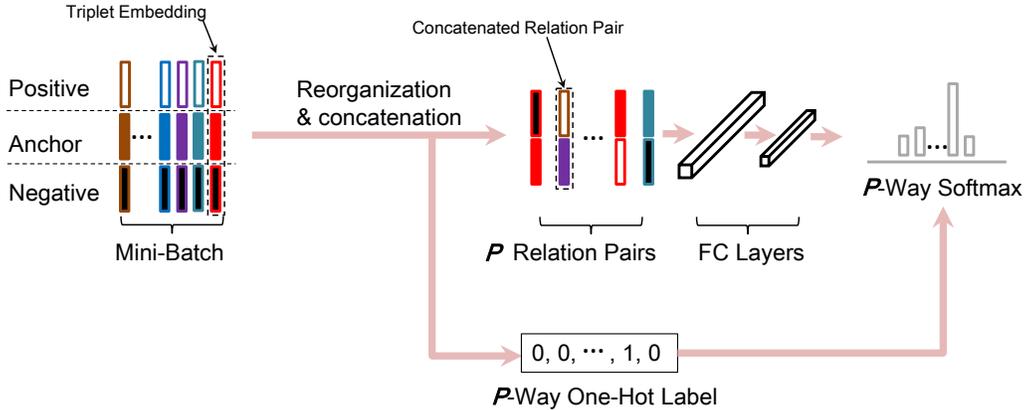


Fig. 6. Detailed illustration of our relation module. Best viewed in color.

### A. Model

FG-SBVR must account for both appearance and motion in fine-grained matching. Therefore, we develop a multi-stream multi-modality joint embedding framework for FG-SBVR that processes appearance and motion using two streams respectively for both modalities. Subsequently, we combine the complementary representations extracted from these two streams to obtain a fused representation. Similarly to several fine-grained retrieval applications [2], [19], we use triplet loss [48] to supervise training. In addition, a softmax cross-entropy loss is added to the relation module sub-network that labels multiple matching and mismatching pairs within one mini-batch. Our model architecture is illustrated in Figure 5, where the four components: Input, Appearance Stream, Motion Stream, and Fusion Mechanism, as well as their relationships are shown. We next detail each of these components.

**Model Input** Within the triplet loss paradigm, training tuples are constructed and each includes: sketch anchor, positive (matching) video, negative (mismatching) video. We explore various triplet construction strategies which will be detailed later. For a given triplet, our operations are: Firstly, we choose one video frame from the positive video as the positive atom of the triplet for appearance stream. Secondly, use this selected frame as a start frame to compute $L$ pairs of consecutive optical flows using the GPU implementation of TV-$L^1$ [49] in OpenCV. Following [37], we calculate the optical flow matrix pair along $x$ and $y$ directions and alternately stack them to form a total of $2L$ input channels. This $2L$ channels of optical flow will work as the positive atom of triplet for motion stream. Thirdly, the same operations are applied to negative atoms for the two streams. Finally, the sketch anchor is separated into "skater" appearance $s^{ap}$ and "motion vector" $s^{mo}$ components, which are then used in the appearance

stream and motion stream as anchors, respectively. Therefore, as shown in Figure 5, there are six input branches in total.

**Appearance Stream CNN**  During training, there are three appearance branches corresponding to the atoms of a triplet $t^{ap} = (s^{ap}, v^{ap,+}, v^{ap,-})$. We assume weight sharing between all three branches (positive and negative frame, appearance sketch). The appearance stream backbone is GoogLeNet Inception V3 [50]. The loss for triplet $t^{ap}$ is:

$$\mathcal{L}_t = \max(0, \triangle + \|\mathcal{F}_{\Theta_{ap,c}}(s^{ap}) - \mathcal{F}_{\Theta_{ap,c}}(v^{ap,+})\|_2^2 \\ - \|\mathcal{F}_{\Theta_{ap,c}}(s^{ap}) - \mathcal{F}_{\Theta_{ap,c}}(v^{ap,-})\|_2^2), \quad (1)$$

where $\triangle$ is a margin and $\mathcal{F}_{\Theta_{ap,c}}$ denotes CNN feature extraction by the appearance stream network, parameterised by $\Theta_{ap,c}$.

**Motion Stream**  Similar to the appearance stream, there are three branches corresponding to triplet atoms $t^{mo} = (s^{mo}, v^{mo,+}, v^{mo,-})$ during training. Due to the different numbers of input channels in sketch and stacked optical flow fields (3 vs. $2L$), we use two CNNs with different input depth. Both are GoogLeNet Inception V3-based, but the latter is modified to use $2L$ input channels. The loss is also a triplet loss combined with a relation loss analogous to Equation 2. We use $\mathcal{F}_{\Theta_{mo,c}}$ to denote CNN feature extraction by the motion stream network, parameterised by $\Theta_{mo,c}$.

**Relation Module**  Limited training sketch-video pairs make a FG-SBVR model vulnerable to overfitting and poor generalisation to test data. Here we introduce our relation module [23], [24] that can be independently applied to both our appearance stream and motion stream to alleviate the data scarcity problem. The idea is that, instead of modelling triplet ranking relationships, we form larger groups and model the more complex group relationship. This aims to maximise the use of the limited training data, because when the group size is larger than 3, there can be many more groups than triplets. More specifically, given a mini-batch, we forward it through our appearance stream or motion stream CNN, and obtain a mini-batch of triplet embedding vectors. Then, we randomly select and reorganise $P$ sketch-video relation pairs, forming one true match pair and $(P-1)$ false match pairs. As shown in Figure 6, we concatenate the embedding vectors for each relation pair, and input it into a relation network consisting of two fully connected layers. We will set the dimensionality of our CNN output embedding as $256D$ in this paper, so that the input dimensionality of our relation module will be $512D$. Simultaneously, the associated ground-truth pairwise relationships are formed as a $P$-dimension one-hot vector as training objective, in which the non-zero element corresponds to the true match pair. We adopt $P$-way cross-entropy softmax loss for our relation module sub-network, denoted as "relation loss" $\mathcal{L}_r$. $P$ is set to 5 in this work. Thus, the total loss for each mini-batch can be defined as

$$\mathcal{L} = \mathcal{L}_t + \lambda_1 \mathcal{L}_r, \quad (2)$$

where $\lambda_1$ is a weighting factor.

In particular, in this work, we independently train a relation module for each stream, parameterised by $\Theta_{ap,r}$ and $\Theta_{mo,r}$ respectively. This is to say that the loss function of each stream

---

**Algorithm 1** Learning algorithm for our proposed multi-stream multi-modality FG-SBVR deep network.

---

**Input:** $D = \{(S_i, V_i)\}_{i=1}^N$.
  1. Train apperence stream as following loop.
  **for** number of training iterations **do**
      1.1 Forward mini-batch through CNN, calculate $\mathcal{L}_t$.
      1.2 Reorganise mini-batch and forward relation pairs through relation module, calculate $\mathcal{L}_r$.
      1.3 Update $\Theta_{ap,c}$ using $\mathcal{L}_t$ and $\mathcal{L}_r$.
      1.4 Update $\Theta_{ap,r}$ using $\mathcal{L}_r$.
  **end for**
  2. Train motion stream as following loop.
  **for** number of training iterations **do**
      2.1 Forward mini-batch through CNN, calculate $\mathcal{L}_t$.
      2.2 Reorganise mini-batch and forward relation pairs through relation module, calculate $\mathcal{L}_r$.
      2.3 Update $\Theta_{mo,c}$ using $\mathcal{L}_t$ and $\mathcal{L}_r$.
      2.4 Update $\Theta_{mo,r}$ using $\mathcal{L}_r$.
  **end for**
**Output:** CNN embedding extractions $\mathcal{F}_{\Theta_{ap,c}}$ and $\mathcal{F}_{\Theta_{mo,c}}$.

---

contains both a triplet term and a relation loss. The detailed training and optimization are described in Algorithm 1.

**Stream Fusion**  Once the triplet ranking and relation module training for the two streams is complete, a natural way to combine them is to concatenate or fuse them with another FC layer, and fine-tune with another triplet loss. However, similarly to the observation in [37], this fails in our case due to overfitting. Thus, we use two fusion approaches to fuse our two streams: (i) ranking-based fusion, and (ii) feature concatenation fusion. For ranking-based fusion, a sketch query $S$ generates a ranked list of matching videos $\{V_j\}$ based on Euclidean distance. We use $r_j^{ap}$ and $r_j^{mo}$ to indicate the rankings of each video $j$ using the appearance and motion stream features respectively. The final ranking $r_j$ of each gallery video clip is the weighted arithmetic mean of its appearance and motion ranks:

$$r_j = \lambda_2 r_j^{ap} + (1 - \lambda_2) r_j^{mo}, \quad (3)$$

where $\lambda_2$ is the weighting factor. For the feature concatenation strategy, we concatenate the features from two streams as the final representation, and then conduct Euclidean distance based ranking.

**Training with Strong Supervision**  Recall that our sketch queries can contain multiple pages corresponding to different segments/sub-clips within the video clip, and that the detailed correspondence is annotated. This provides the strongest supervision for our task. Specifically, for each single sketch anchor, its positive video candidates are frames within the corresponding sub-clip. Frames outside the corresponding sub-clip, or frames in different clips entirely, are treated as negative.

*B. Weakly Supervised Learning*

The conventional strongly supervised setting requires labor intensive cross-modal annotation (Figure 2). It would be easier
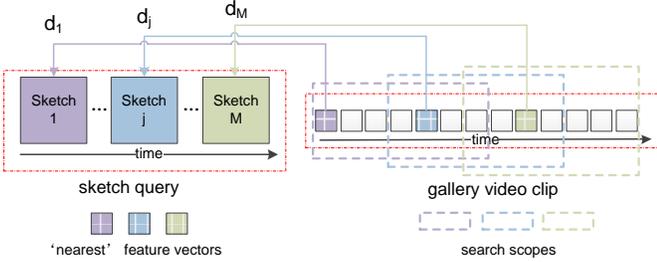
Fig. 7. Illustration for the similarity distance calculation for sketch query and gallery video. Best viewed in color.

work, sketches and videos are respectively represented as: $S = \{s_j\}_{j=1}^{M} = \{(F_{\Theta_{ap,c}}(s_j^{ap}), F_{\Theta_{mo,c}}(s_j^{mo}))\}_{j=1}^{M}$ and $V = \{v_k\}_{k=1}^{O} = \{(F_{\Theta_{ap,c}}(v_k^{ap}), F_{\Theta_{mo,c}}(v_k^{mo}))\}_{k=1}^{O}$ respectively. A simple solution is to choose the match that has the lowest sum of nearest neighbour matching costs:

$$D(S,V) = \frac{1}{M} \sum_{j=1}^{M} \min_{k \in [\varphi(j), \psi(j)]} d(s_j, v_k),$$

$$\left( \begin{cases} \varphi(j) = 1 \ and \ \psi(j) = O/2, & j = 1 \\ \varphi(j) = O/2 \ and \ \psi(j) = O, & j = M \\ \varphi(j) = O/4 \ and \ \psi(j) = 3O/4, \ other \end{cases} \right),$$

(4)

where $\varphi(j)$ and $\psi(j)$ are lower and upper bounds of $k$. Figure 7 provides an illustration.

## V. EXPERIMENTS

### A. Experiment Settings

**Dataset Split**  We generate training, validation, and testing sets by randomly splitting the 528 clips into 350 (with 971 sketches) for training, 50 (with 131 sketches) for validation, and 128 (with 346 sketches) for testing. Thus during testing, we have 128 sketch sequences and 128 video clips as the queries and gallery, respectively.

**Implementation Details**  All experiments are implemented in PyTorch, and run on a single TITAN Xp GPU. We use model hyperpameters values obtained using 5-fold cross-validation on the training set: $\delta = 0.5, L = 5, T\% = 0.1, \lambda_1 = 0.001, \lambda_2 = 0.5$. RMSprop optimizer is used with initial learning rate 0.001 and mini-batch size is 16. We initialise our branches with Inception V3 with the following modifications: The original "fc" components of Inception are replaced by two new fully-connected layers ($2048D \rightarrow 512D$, $512D \rightarrow 256D$). For the motion stream, the first layer with non-standard depth is randomly initialised. Our relation module sub-network has two fully-connected layers ($512D \rightarrow 128D$, $128D \rightarrow 32D$). The input size of the appearance stream and sketch motion sub-stream is $3 \times 299 \times 299$. The optical flow sub-stream input is $10 \times 299 \times 299$.

**Evaluation Metric**  Given our instance-level retrieval task, we use retrieval accuracy as a metric. We quantify this by cumulative matching accuracy at various ranks. Acc@K is the proportion of sketch sequence queries whose true-match video clips are ranked in the top K (1 means 100%).

**Competitors**  FG-SBVR is a brand new problem, thus no existing methods can be compared directly here. The main competitors here are the state-of-the-art video analysis models extended for the FG-SBVR problem. More specifically, for fair comparison, we still employ a multi-stream multi-modality network. But both streams of the video modality are now based on the latest 3D convolutional networks including Non-Local Neural Network [20], 3D ResNet18 [21], and T-C3D [22]. The two sketch streams are unchanged and aligned with the duplicated video streams in the hope that the static and dynamic aspects of the video contents can be disentangled through the alignment process, trained with the same triplet and relation module losses. In addition, for an ablation study, we also compare several variants of our own model. In particular, we compare our model trained with different types

to scale FG-SBVR if the model could be learned with less detailed and intensive annotation. Weakly supervised learning (WSL) for FG-SBVR is thus required. This can be formulated as a multi-instance learning (MIL) [51] problem. Given an anchor sketch-page, there will be video clips to which it definitely does not correspond (negative bag), and those to which it will correspond to at least one frame within the clip (positive bag). Each of those clips (bags) contains multiple frames (instances). In particular, we consider that we have the sketch-sequence to video-clip pairing, but not the detailed page-frame level pairing (Figure 2, green lines). In this case, there is one positive bag per anchor page, *i.e.*, the corresponding clip. All non-corresponding clips are negative bags. The challenge is to correctly estimate the positive/matching, and negative/mismatching instances (frames) within each positive bag (video).

**Pruning Heuristics and Initialisation**  To prune the search space of positive bags defined above, we further consider: (i) If a query sketch is the first in a sequence, its positive bag consists of the first 50% of frames in a potentially matching video according to the criteria above. (ii) If a query sketch is last in a sequence, its positive bag consists of the last 50% frames in a potentially matching video. (iii) Otherwise the positive bag is the middle 50% of clip frames. Both appearance and motion streams share the same criteria. All instances (frames) within positive bags are initialised as positive.

**Multi-Instance Learning**  Given the above bag definition and initialisation, we iteratively refine the labels of instances within positive bags using multi-instance training. As per conventional MIL [51], training alternates between phases of classifier/representation learning, and re-estimation of positive instances as follows: (i) Learning: Update the network with back-propagation and loss $\mathcal{L}$, assuming the positive/negative instance labels are fixed. (ii) Label Inference: Set the least likely matches (furthest $T\%$ distance) frames in positive bags to be negative, assuming the network parameters are fixed.

### C. Model Deployment

Given a trained model, during testing we need to match a sketch-sequence to a video clip (frame-sequence). However, our network has learned to rank sketch-pages and video-frames. Thus, we need to define how to aggregate a set of pairwise page-frame scores. In testing, we use the CNN output embedding to conduct retrieval. In particular, after deep feature extraction by our two-stream net-

TABLE II
ABLATIVE EVALUATION OF OUR FG-SBVR MODEL UNDER STRONG SUPERVISION. CHANCE LEVEL PERFORMANCE IS 0.0078 ($\approx 1/128$) ACC.@1.

| Model | Triplet Ranking | | | Triplet Ranking + Relation Module | | |
|---|---|---|---|---|---|---|
| | acc.@1 | acc.@5 | acc.@10 | acc.@1 | acc.@5 | acc.@10 |
| app. stream | 0.1250 | 0.3281 | 0.4063 | 0.1719 | 0.3516 | 0.5156 |
| motion stream | 0.1406 | 0.3438 | 0.5469 | 0.1719 | 0.3828 | 0.5781 |
| ranking fusion | 0.2188 | 0.4141 | 0.5547 | 0.2969 | 0.6094 | 0.7344 |
| concat fusion | 0.3047 | 0.6172 | 0.7344 | 0.3438 | 0.6094 | 0.7656 |

TABLE III
FG-SBVR RETRIEVAL RESULTS OBTAINED WITH WEAK SUPERVISION.

| Model | Triplet Ranking | | | Triplet Ranking + Relation Module | | |
|---|---|---|---|---|---|---|
| | acc.@1 | acc.@5 | acc.@10 | acc.@1 | acc.@5 | acc.@10 |
| app. stream | 0.0234 | 0.1016 | 0.1719 | 0.0469 | 0.1094 | 0.1641 |
| motion stream | 0.0469 | 0.1250 | 0.2500 | 0.0703 | 0.1797 | 0.2656 |
| ranking fusion | 0.0547 | 0.1094 | 0.1641 | 0.0859 | 0.1172 | 0.1719 |
| concat fusion | 0.0625 | 0.1094 | 0.1953 | 0.0703 | 0.1250 | 0.1875 |

TABLE IV
FG-SBVR RESULTS OF 3D CNN BASED BASELINES UNDER STRONG
SUPERVISION SETTING.

| 3D CNN | Model | acc.@1 | acc.@5 | acc.@10 |
|---|---|---|---|---|
| | app. stream | 0.0938 | 0.2109 | 0.3359 |
| Non-Local [20] | motion stream | 0.0703 | 0.1953 | 0.3047 |
| | ranking fusion | 0.1016 | 0.1875 | 0.3281 |
| | concat fusion | 0.1016 | 0.2969 | 0.4141 |
| | app. stream | 0.0469 | 0.0781 | 0.1094 |
| 3D ResNet18 [21] | motion stream | 0.0469 | 0.1016 | 0.1563 |
| | ranking fusion | 0.0234 | 0.0781 | 0.1250 |
| | concat fusion | 0.0547 | 0.0859 | 0.1328 |
| | app. stream | 0.0391 | 0.1094 | 0.1719 |
| T-C3D [22] | motion stream | 0.0234 | 0.0938 | 0.1641 |
| | ranking fusion | 0.0313 | 0.0859 | 0.1484 |
| | concat fusion | 0.0313 | 0.0703 | 0.1719 |

TABLE V
PERFORMANCE IMPROVEMENT ON NON-LOCAL NEURAL NETWORK [20]
ACHIEVED BY INTRODUCING THE RELATION MODULE.

| Model | acc.@1 | acc.@5 | acc.@10 |
|---|---|---|---|
| app. stream | 0.0235 | 0.0781 | 0.0703 |
| motion stream | 0.0156 | 0.0234 | 0.0313 |
| ranking fusion | 0.0078 | 0.0235 | 0.0000 |
| concat fusion | 0.0157 | 0.0625 | 0.0313 |

of supervision, including strong supervision (see Section IV-A) and weak supervision (see Section IV-B). Since our model has static appearance and dynamic motion streams, we also compare variants of our model with one stream only and the full model with the two-streams fused in different ways (see Section IV-A). Note that we also attempted to put RNN on top of the video and sketch streams to explicitly encode the temporal order information following [36], but the results are much worse so not reported here.

## B. Results

**Strong Supervision** Under strong supervision, we first conduct an ablative study on the contributions of the multi-stream multi-modal architecture design and relation module in our model. The following observations can be made from Table II: (i) In terms of single-stream performance, motion outperforms appearance particularly at higher ranks. This is interesting: Each motion vector sketch component contains only 1 or 2 strokes, whilst the mean number of strokes for the skater sketch component is around 100 (see Table I). Those 1-2 motion strokes seem worth 100 strokes that are used for depicting static appearance of the skater and her representative posture. (ii) When the two streams are fused using either ranking or feature concatenation, the performance is improved significantly. These results confirm that appearance and motion patterns contain complementary information for FG-SBVR. Further, among the two fusion strategies, feature concatenation is clearly more effective. (iii) When the relation module is added, large improvements on retrieval accuracy are obtained.

Next, we compare our full model with a number of baselines extended from state-of-the-art 3D CNN models. Table IV shows the results obtained by three models based on 3D CNNs introduced in 2018 deployed in the same multi-stream multi-modality network, trained with the same triplet ranking and relation module. Comparing Table IV with Table II, it is clear that our model with the video modality decomposed explicitly into dynamic and static parts, modelled with optical flow and image CNN respectively, is much better (24% higher on acc.@1). These results thus further validate our model design. It is also interesting to note that for non-local network [20], stream fusion helps and concatenation based fusion is the most effective strategy.

**Relation Module Analysis** To understand the impact of our relation module, we also study its efficacy in combination with the baselines. Specifically, we pick the strongest baseline, Non-Local network [20] and compute the performance difference

TABLE VI
IMPACT OF NUMBER OF PAIRS PER MINI-BATCH USED BY RELATION MODULE.

| model | Num paris per mini-batch | acc.@1 | acc.@5 | acc.@10 |
|---|---|---|---|---|
| strong supervision app. stream (SAS) | 5 | 0.1719 | 0.3516 | 0.5156 |
| strong supervision app. stream (SAS) | 10 | 0.1953 | 0.3672 | 0.5234 |
| strong supervision app. stream (SAS) | 15 | 0.2031 | 0.4531 | 0.5859 |
| strong supervision app. stream (SAS) | 20 | 0.2344 | 0.5156 | 0.6094 |

TABLE VII
SKETCH-BASED ACTION DETECTION RESULTS. CHANCE PERFORMANCE
(ACC.@1) IS 0.0029 ($\approx 1/346$).

| Supervision | Model | Accuracy |
|---|---|---|
| Strong Supervision | app. stream | 0.4133 |
| | motion stream | 0.3382 |
| | ranking fusion | 0.4682 |
| | concat fusion | 0.4798 |
| Weak Supervision | app. stream | 0.1185 |
| | motion stream | 0.1936 |
| | ranking fusion | 0.1850 |
| | concat fusion | 0.1647 |

between the models with and without the relation module. Table V shows that in most cases, the relation modulebrings improvement, indicating the general applicability of such meta-learning inspired techniques for dealing with scarce training data.

One reason for the efficacy of our relation module is that it leverages more negative pairs within each mini-batch. Table VI analyses the relation module from this perspective, showing that increasing the relation pairs per mini-batch leads to improved retrieval performance. For our main results in this paper, we use 5 relation pairs within each mini-batch due to limited GPU memory.

**Weak Supervision** Table III also shows the results obtained when our model is trained with weak supervision. It is clear that: (i) Like the strong supervision setting, single motion stream outperforms single appearance stream. (ii) Two-stream fusion now does not guarantee to improve performance, due to the poor performance of the appearance stream. (iii) All the accuracy values are significantly lower than that under strong supervision, as expected. Overall, it is worth pointing out that there is a large scope for further improvement under this challenging setting, in order to narrow the distinct performance gap between supervised and weakly supervised models.

**Fine-Grained Sketch-Based Action Detection** Going beyond FG-SBVR, we can also use our method and dataset to explore an even more fine-grained application, namely fine-grained instance-level sketch-based action detection (FG-SBAD). Given a motion sketch-page and a video clip, the goal of FG-SBAD is to localise the target action depicted by the sketch. We propose a straightforward solution to FG-SBAD: traverse video clip frame-by-frame and report the index of the nearest-neighbour frame. If the proposal is in the range

(*i.e.*, within 5 frames) of the sketch-to-frames ground truth, then it is regarded as a successful detection. As illustrated in Table VII, the strongly-supervised concatenation fusion approach performs much better than the weakly supervised alternative.

**Qualitative Results** We next show some visual examples of the retrieval results obtained using our multi-stream multi-modality model and its variants. In Figure 8, a sequence of three sketches are used as query and the top-10 ranked videos using different models are shown. The query sketch sequence captures the key moments of the video sequence. In particular, the motion vector parts of the three sketches indicate one spin movement and two glissade movements, respectively; in the meantime, the skater parts of the sketches contain visual details of the skater's appearance such as stripes on the clothes and glove, as well as her body posture at those key moments. From the retrieval results, it can be seen that: (i) The multi-stream fusion models (feature concatenation based fusion and ranking based fusion) give the desired results – the true match is ranked at the top. (ii) From the retrieval results of the appearance-stream model variant, we can see that although the correct match is not in the top 10, all skaters in the top 10 retrieved videos wear gloves or single shoulder dress, similar to those of the skater in the true match video. These results suggest that without the motion vector, the skater part of the sketches is not discriminative enough for the model to find the correct skating sequence – the model put too much emphasise on the static appearance of the skater rather than her movements. (iii) By contrast, the motion-stream model variant is able to retrieval video sequences containing similar spin or glissade movements with the given query. However, without any information about the static appearance of the skater, the model is unable to distinguish video sequences of similar skating routines but performed by different skaters.

**Localisation** Note that our two-stream model is able to produce a similarity/matching score between a query sketch and each frame of a video sequence (*i.e.*, Fine-Grained Sketch-Based Action Detection as discussed in Section 5 of the main paper). In Figure 9, we show that given a query sketch, which frame in the correctly matched video sequence has the highest matching score. The results suggest that, if the video sequence can be correctly retrieved, our model can be used to accurately localize which specific time of the sequence the query sketch is depicting. In particular, the body pose of the skater in the best matched frame is remarkably similar to the body pose of the sketched skater.

Fig. 8. Qualitative comparison of top 10 retrieval results using different variants of our models under the strongly supervised setting. The 10 videos are ordered from top to bottom and from left to right according to their ranks. The true matches are highlighted in green.

Fig. 9. Results for sketch-based action detection of our concatenation fused full model under the strongly supervised setting. In each row, sketch is on the left, and the green-bordered frame is best matched frame. The neighbouring six frames are also shown.

**Running Cost** All our experiments are conducted on an Intel Core i7-7700K 4.2 GHz CPU and single TITAN Xp GPU. Training our two-stream model with strong supervision and weak supervision takes about 20 hours and 30 hours, respectively. For testing, it takes about 200 milliseconds for one retrieval.

## VI. CONCLUSION AND FUTURE WORK

We introduced the novel task of fine-grained instance-level sketch-based video retrieval (FG-SBVR) and a dataset to enable the research on this task. We also proposed a novel multi-stream multi-modality network with relation module to solve this problem in both strongly- and weakly-supervised settings. Our new dataset can also support future research on tasks such as video summarisation, sketch-based video generation, and multi-modal tasks that combine visual and audio cues (via commentary track). For example, the dataset can be used directly for video summarisation to complement existing datasets such as TVSum [52] and CoSum [53]. The recent research on video-to-video synthesis [54] can now be extended to cross-modal video-to-sketch and sketch-to-video synthesis using our dataset.

## REFERENCES

[1] P. Xu, "Deep learning for free-hand sketch: A survey," *arXiv preprint arXiv:2001.02600*, 2020.

[2] Q. Yu, F. Liu, Y.-Z. Song, T. Xiang, T. M. Hospedales, and C.-C. Loy, "Sketch me that shoe," in *CVPR*, 2016.

[3] P. Xu, Q. Yin, Y. Huang, Y.-Z. Song, Z. Ma, L. Wang, T. Xiang, W. B. Kleijn, and J. Guo, "Cross-modal subspace learning for fine-grained sketch-based image retrieval," *Neurocomputing*, 2018.

[4] J. Collomosse, T. Bui, and H. Jin, "Livesketch: Query perturbations for guided sketch-based visual search," in *CVPR*, 2019.

[5] D. Ha and D. Eck, "A neural representation of sketch drawings," in *ICLR*, 2018.

[6] R. G. Schneider and T. Tuytelaars, "Example-based sketch segmentation and labeling using crfs," *ACM TOG*, 2016.

[7] P. Xu, Y. Huang, T. Yuan, K. Pang, Y.-Z. Song, T. Xiang, T. M. Hospedales, Z. Ma, and J. Guo, "Sketchmate: Deep hashing for million-scale human sketch retrieval," in *CVPR*, 2018.

[8] U. R. Muhammad, Y. Yang, Y.-Z. Song, T. Xiang, and T. M. Hospedales, "Learning deep sketch abstraction," in *CVPR*, 2018.

[9] Y. Ye, Y. Lu, and H. Jiang, "Human's scene sketch understanding," in *ICMR*, 2016.

[10] Y. Xie, P. Xu, and Z. Ma, "Deep zero-shot learning for scene sketch," *arXiv preprint arXiv:1905.04510*, 2019.

[11] P. Xu, Z. Song, Q. Yin, Y.-Z. Song, and L. Wang, "Deep self-supervised representation learning for free-hand sketch," *arXiv preprint arXiv:2002.00867*, 2020.

[12] E. Tulving and D. Murray, "Elements of episodic memory," *Canadian Psychology*, 1985.

[13] K. P. Madore, H. G. Jing, and D. L. Schacter, "Selective effects of specificity inductions on episodic details: evidence for an event construction account," *Memory*, 2018.

[14] J. P. Collomosse, G. McNeill, and Y. Qian, "Storyboard sketches for content based video retrieval," in *ICCV*, 2009.

[15] R. Hu and J. Collomosse, "Motion-sketch based video retrieval using a trellis levenshtein distance," in *ICPR*, 2010.

[16] R. Hu, S. James, and J. Collomosse, "Annotated free-hand sketches for video retrieval using object semantics and motion," in *MMM*, 2012.

[17] R. Hu, S. James, T. Wang, and J. Collomosse, "Markov random fields for sketch based video retrieval," in *ICMR*, 2013.

[18] S. James and J. Collomosse, "Interactive video asset retrieval using sketched queries," in *European Conference on Visual Media Production*, 2014.

[19] P. Sangkloy, N. Burnell, C. Ham, and J. Hays, "The sketchy database: learning to retrieve badly drawn bunnies," *ACM TOG*, 2016.

[20] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *CVPR*, 2018.

[21] K. Hara, H. Kataoka, and Y. Satoh, "Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet?" in *CVPR*, 2018.

[22] K. Liu, W. Liu, C. Gan, M. Tan, and H. Ma, "T-c3d: Temporal convolutional 3d network for real-time action recognition," in *AAAI*, 2018.

[23] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, "Learning to compare: Relation network for few-shot learning," in *CVPR*, 2018.

[24] W. Xie, L. Shen, and A. Zisserman, "Comparator networks," in *ECCV*, 2018.

[25] Y.-P. Tan, S. R. Kulkarni, and P. J. Ramadge, "A framework for measuring video similarity and its application to video query by example," in *ICIP*, 1999.

[26] Y. Li, R. Wang, Z. Huang, S. Shan, and X. Chen, "Face video retrieval with image query via hashing across euclidean space and riemannian manifold," in *CVPR*, 2015.

[27] A. Araujo and B. Girod, "Large-scale video retrieval using image queries," *TCSVT*, 2018.

[28] R. Xu, C. Xiong, W. Chen, and J. J. Corso, "Jointly modeling deep video and compositional text to bridge vision and language in a unified framework." in *AAAI*, 2015.

[29] N. C. Mithun, J. Li, F. Metze, and A. K. Roy-Chowdhury, "Learning joint embedding with multimodal cues for cross-modal video-text retrieval," in *ICMR*, 2018.

[30] B. Coskun, B. Sankur, and N. Memon, "Spatio–temporal transform based video hashing," *TMM*, 2006.

[31] M. Li and V. Monga, "Robust video hashing via multilinear subspace projections," *TIP*, 2012.

[32] G. Ye, D. Liu, J. Wang, and S.-F. Chang, "Large-scale video hashing via structure learning," in *ICCV*, 2013.

[33] Z. Chen, J. Lu, J. Feng, and J. Zhou, "Nonlinear structural hashing for scalable video search," *TCSVT*, 2018.

[34] J. Song, H. Zhang, X. Li, L. Gao, M. Wang, and R. Hong, "Self-supervised video hashing with hierarchical binary auto-encoder," *TIP*, 2018.

[35] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venu-gopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *CVPR*, 2015.

[36] B. Singh, T. Marks, M. Jones, O. Tuzel, and M. Shao, "A multi-stream bi-directional recurrent neural network for fine-grained action detection," in *CVPR*, 2016.

[37] K. Simonyan and A. Zisserman, "Two-stream convolutional networks for action recognition in videos," in *NIPS*, 2014.

[38] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool, "Temporal segment networks: Towards good practices for deep action recognition," in *ECCV*, 2016.

[39] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *ICCV*, 2015.

[40] Z. Qiu, T. Yao, and T. Mei, "Learning spatio-temporal representation with pseudo-3d residual networks," in *ICCV*, 2017.

[41] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *CVPR*, 2017.

[42] L. Wang, W. Li, W. Li, and L. Van Gool, "Appearance-and-relation networks for video classification," in *CVPR*, 2018.

[43] A. Diba, M. Fayyaz, V. Sharma, M. Mahdi Arzani, R. Yousefzadeh, J. Gall, and L. Van Gool, "Spatio-temporal channel correlation networks for action classification," in *ECCV*, 2018.

[44] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy, "Rethinking spatiotem-poral feature learning: Speed-accuracy trade-offs in video classification," in *ECCV*, 2018.

[45] Y. Chen, Y. Kalantidis, J. Li, S. Yan, and J. Feng, "Multi-fiber networks for video recognition," in *ECCV*, 2018.

[46] M. Eitz, J. Hays, and M. Alexa, "How do humans sketch objects?" *ACM TOG*, 2012.

[47] G. A. Miller, "The magical number seven, plus or minus two: Some limits on our capacity for processing information." *Psychological review*, 1956.

[48] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embed-ding for face recognition and clustering," in *CVPR*, 2015.

[49] C. Zach, T. Pock, and H. Bischof, "A duality based approach for realtime $tv - l^1$ optical flow," in *Joint Pattern Recognition Symposium*, 2007.

[50] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *CVPR*, 2016.

[51] S. Andrews, I. Tsochantaridis, and T. Hofmann, "Support vector ma-chines for multiple-instance learning," in *NIPS*, 2003.

[52] Y. Song, J. Vallmitjana, A. Stent, and A. Jaimes, "Tvsum: Summarizing web videos using titles," in *CVPR*, 2015.

[53] W.-S. Chu, Y. Song, and A. Jaimes, "Video co-summarization: Video summarization by visual co-occurrence," in *CVPR*, 2015.

[54] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro, "Video-to-video synthesis," in *NeurIPS*, 2018.