

# What and How Well You Performed? A Multitask Learning Approach to Action Quality Assessment

Paritosh Parmar

Brendan Tran Morris

University of Nevada, Las Vegas

parmap1@unlv.nevada.edu, brendan.morris@unlv.edu

## Abstract

Can performance on the task of action quality assessment (AQA) be improved by exploiting a description of the action and its quality? Current AQA and skills assessment approaches propose to learn features that serve only one task - estimating the final score. *In this paper, we propose to learn spatio-temporal features that explain three related tasks - fine-grained action recognition, commentary generation, and estimating the AQA score.* A new multitask-AQA dataset, the largest to date, comprising of 1412 diving samples was collected to evaluate our approach (<http://rtis.oit.unlv.edu/datasets.html>). We show that our MTL approach outperforms STL approach using two different kinds of architectures: C3D-AVG and MSCADC. The C3D-AVG-MTL approach achieves the new state-of-the-art performance with a rank correlation of 90.44%. Detailed experiments were performed to show that MTL offers better generalization than STL, and representations from action recognition models are not sufficient for the AQA task and instead should be learned.

## 1. Introduction

What score should an athlete receive on her dive/gymvault/skating/etc? Which med student has the highest surgical skill level? How well can he paint or draw? How is a patient progressing in their physical rehabilitation program? Answering these questions involves the quantification of the quality of the action – determining *how well* the action was carried out, also known as action quality assessment (AQA). Existing AQA [18, 16, 26, 13, 25] and skills assessment [4, 10, 31, 32, 33] approaches use a single label, known as a final score or skill-level, to train the system using some kind of regression or ranking loss function. However, the performance of these systems is limited and it seems that a single score is not sufficient to characterize a complicated action. In AQA, the final

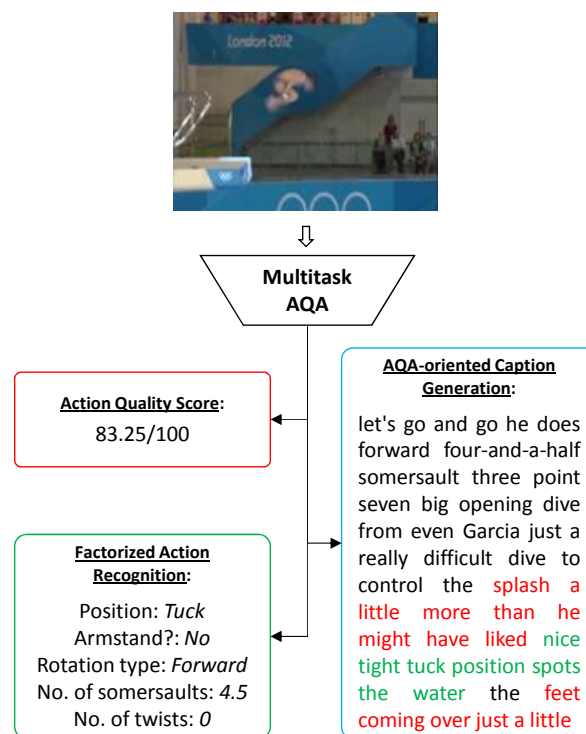


Figure 1: **Multitask AQA concept.** Recognizing an action instance in detail and verbally describing its good and bad points can be helpful in the process of quantifying the quality of that action instance. We propose to learn a model that delineates an action besides measuring its quality. *To see the videos play, please download the manuscript and view in an Adobe Reader.*

score is dependent on what was done (this determines the difficulty level) and how was that done (this determines the quality of execution). We pose the following question: *can learning to describe and commentate on the action instances help improve the performance on the AQA task?*

We hypothesize that by forcing the network to learn to do so will help better characterize the action, and hence aid

in AQA. So, rather than using just a single encompassing quality label to train the network, we introduce a multitask learning (MTL) approach (Fig. 1) to assess the quality of an action. Specifically, we propose to utilize 3D CNN’s to learn spatio-temporal representations of salient motion and appearance; optimize those using loss functions which account for i) the action quality score, ii) factorized (detailed) action classification, and iii) generate a verbal commentary of performance; and are trained end-to-end. **Note that the architectures are multitask and not multi-modal since the input does not use captions or action classification to produce the AQA score.** Besides straight forward utility for AQA and action classification, automatic commentary or sports narrative generation has been viewed valuable and greatly applicable in a recent work by Yu *et al.* [29].

For AQA tasks, domain experts can provide detailed analysis of performance. In the professional sports setting, ground truth annotations for detailed action classification and commentary by former athletes are readily available in broadcast footage facilitating extraction of labels and descriptive captions. As such, to evaluate our approach, we introduce the first multitask AQA dataset with 1412 samples of diving which is also the largest AQA dataset to date.

Experimental evaluation show that performance of both the architectures improved as more tasks were added and the C3D-AVG-MTL variant outperforms all existing AQA approaches in literature. MTL was shown to outperform STL across various training set sizes. Further experiments explore the AQA-orientedness of the feature representations learned by our networks and find they outperform action-recognition representations on unseen actions indicating that better generalized concepts of quality were learned.

**Contributions:** primary novelty of this works lies in the problem formulation – to learn spatio-temporal representations by optimizing networks end-to-end jointly for fine-grained action description and AQA scoring. Task selection is intuitive. No previous work has done this; not just for AQA, but even action recognition and captioning tasks. We release a novel MTL-AQA dataset which is the largest AQA dataset so far, much more diverse, challenging, and richly annotated with factorized fine-grained action class and AQA-oriented captions. Our dataset can help researchers in the field to examine new ideas for AQA and auxiliary tasks. We show that our MTL approach works across different architectures. Our approach is applicable to a wide range of problems. Our proposed models are simple, yet intuitive, and effective in carrying out central of learning representations in a MTL setting by optimizing networks end-to-end. **Our C3D-AVG-MTL surpasses all the existing approaches.**

## 2. Related Work

**AQA:** Pirsiavash *et al.* [18] proposed the use of DFT/DCT of body pose as features for a support vector regressor (SVR) to map to a final action quality score. They introduced an action quality dataset containing two actions: Diving and Figure Skating. However, since their method relied solely on pose features, it neglected important visual quality cues, like splash in the case of Diving. Since accurate pose is especially difficult in sports scenarios where athletes undergo extremely convoluted poses, Venkataraman *et al.* [25] better encoded using the approximate entropy of the poses to improve the results.

More recently, spatio-temporal features from 3D convolutional neural networks (C3D) [24] proved to be very successful on a related task of action recognition since they captured appearance and salient motion. Seeing this as a desirable property that would help to take into account visual cues, Parmar and Morris [16] proposed using C3D features for AQA. They proposed three frameworks, C3D-SVR, C3D-LSTM, and C3D-LSTM-SVR, which differed in their feature aggregation and regression scheme. All the frameworks worked better than previous models proving the efficacy of C3D features for AQA. Xiang *et al.* [26] proposed breaking video clips into action specific segments and fusing segment-averaged features instead of over full videos. By adding finer segment labels to data samples performance was improved. Li *et al.* [13] divide a sample into 9 clips and use 9 different C3D networks dedicated to different stages of Diving. Features are concatenated and further processed through `conv` and `fc` layers to produce a final AQA score using a ranking loss along with the more typical L2 loss. Xu *et al.* [27] tackle AQA for longer action sequences using self-attentive and multiscale convolutional skip LSTM.

**Skills assessment:** Zia *et al.* [33] extract spatio-temporal interest points (STIP’s) in the frequency domain to classify a sample into novice, intermediate or expert skills level. Instead of using handcrafted STIP’s Doughty *et al.* [4] learn and use convolutional features with ranking loss as their objective function to evaluate surgical, drawing, chopstick use and dough rolling skills. In their subsequent work [5], they use temporal attention. Li *et al.* [14], make use of spatial attention in the assessment of hand manipulation skills. Bertasius *et al.* [1] focus on measuring basketball skills but rely only on assessment of a single basketball coach making their dataset subjective to a particular evaluator.

All of the existing AQA and SA frameworks are single task models and only give the final AQA score. Our proposed framework is a multitask model to recognize the action, measures its quality and also generates captions (or

Dataset	Events	Height	Genders	# Samples	Events	View Variation/ Background	Labels
MIT Dive [18]	Individual	10m Platform	Male	159	1	No/Same	AQA score
UNLV Dive [16]	Individual	10m Platform	Male	370	1	No/Same	AQA score
Ours MTL-AQA	Individual, Synchronous	3m Springboard, 10m Platform	Male, Female	1412	16	Yes/Different	AQA score, Action class, Commentary

Table 1: **Details of our newly introduced dataset**, and its comparison with the existing AQA datasets.

Position	Armstand	Rotation type	# SS	# TW
Free	No	Inward	0 to 4.5	0 to 3.5
Tuck	Yes	Reverse		
Pike		Backward		
		Forward		

Table 2: **Classification of dives**. Each combination of the presented sub-fields produces a different kind of maneuver.

commentary).

**Multi-modal approaches and captioning:** Images and videos (especially sports) are often accompanied by a caption or commentary which can themselves serve as labels yet to be exploited for AQA or skill assessment. Quattoni *et al.* [19] use large quantities of unlabeled images, with associated captions, to learn image representations. They found that this sort of pre-training with extra information could speed up the learning on a target task. Rather than using captions as groundtruth labels, Sonal *et al.* [6] treated captions as a “view” and use them along with images to learn a classifier using co-training. They again used commentary as a “view” for action recognition with success. To train an activity classifier in an automated fashion, without the requirement of any manual labeling, Sonal and Mooney [7] make use of broadcast closed captions and used the system for video retrieval. There are a few works which focus on captioning in sports settings. Yu *et al.* [29] address the task of generating fine-grained video descriptions for basketball and evaluate performance using their novel metric. Commentary generation in cricket has been addressed in [20, 21], while Sukhwani addressed the problem of describing tennis videos in [23]. While these works focus on captioning or improving captioning, we integrate a captioning task with an AQA task to provide stronger supervision as commentary is a verbal description of AQA.

### 3. Multitask AQA Dataset

In order to facilitate research in the area of AQA, we release a new dataset. This is the first of a kind multitask

AQA dataset. With 1412 samples, it is the largest AQA dataset to date. This particular dataset focuses only on Diving as it has seen the most usage recently. Data was compiled from 16 different events unlike the single main event (2012 Olympics Men’s 10m Platform Diving competition) used for previous datasets [18, 16] to provide significantly more variation. Diving samples in the new dataset were collected from various International competitions and include the 10m Platform as well as 3m Springboard, include both male and female athletes, individual or pairs of synchronized divers, and different views. A comparison of our new dataset with existing Diving AQA sets is provided in Table 1.

Since data was collected from televised international events, before the athletes perform their routines, information regarding their routine is displayed. This information includes the difficulty of the dive and a description of the dive. The AQA score is extracted from the judges’ scores after the dive completion. The dataset uses the same dive classification strategy as Nibali *et al.* [15], where instead of using dive number (equivalent to an action class in action recognition) directly, we factorize a dive into its components such as the position of the dive, the number of somersaults (SS), and number of twists (TW). Full details for the dive classification is in Table 2.

Further, during and after a diving routine, television analysts provide commentary. These analysts are often retired athletes and have deep understanding of the sport. This verbal account of the athlete’s performance is recorded for the third type of action label. The commentary was considered an important indicator for performance since it was the only way to “watch” an event before telecast was available. Commentators say what the athlete performed, what was correct with the athlete’s performance, and where and how athletes made mistakes. This provides deeper insight into the athlete’s performance and can help an average person better understand the sport. We used Google’s Speech-To-Text API to convert commentary audio to text.

## 4. Multitask Approach to AQA

MTL is a machine learning paradigm in which a single model caters to more than a single task. An example is to recognize road signs, roads, and vehicles together while an STL approach would require separate models for each object type. MTL tasks are generally chosen such that they are related to one another and their networks have a common body that branches into task-specific heads. The total network loss is the sum of individual task losses. When optimized end-to-end, the network is able to learn richer representation in the common body section since it must be able to serve/explain all tasks. With the use of related auxiliary tasks, which are complementary to the main task, the richer representation tends to help improve performance on the main task.

In general, not just for diving, action quality is a function of *what* action was carried out and *how well* that action was executed. This makes the choice of auxiliary tasks natural: detailed action recognition is the answer to the ‘*what*’ part and commentary, being a verbal description containing good and bad points about action execution, is an answer to the ‘*how well*’ part. AQA can be thought of as finding a function that maps input video to the AQA scores. Caruana in [2] views supervision signals from auxiliary tasks as an inductive bias (assumptions). Inductive bias can be thought of as constraints that restrict the hypothesis/search space when finding the AQA function. Through inductive biases, MTL provides improved generalization as compared STL [2].

In this work, the main task is to assess the action quality (AQA score) and the auxiliary tasks are to recognize the action (dive type classification) and to generate descriptive captions/commentary. Action recognition in turn consists of five fine-grained dive sub-recognition tasks: recognizing position and rotation type, detecting armstand, and counting somersaults and twists.

First, let us formalize the settings and objective functions. AQA is a regression problem where, generally, the Euclidean distance between the predicted quality score and the ground truth is used as the objective function to be minimized [16, 26, 13]. Initial experimentation found that using L1 distance in addition to L2 yielded better results on the AQA task

$$\mathcal{L}_{AQA} = -\frac{1}{N} \sum_{i=1}^N (x_i - y_i)^2 + |x_i - y_i| \quad (1)$$

where  $x_i$  is the predicted score and  $y_i$  is the ground truth score for each of the  $N$  samples. For action recognition, we use cross-entropy loss between the predicted labels and

ground truth label

$$\mathcal{L}_{Cls} = -\frac{1}{N} \sum_{i=1}^N \sum_{sa} \sum_{j=1}^{k_{sa}} y_{i,j}^{sa} \log(x_{i,j}^{sa}) \quad (2)$$

where  $k_{sa}$  is the number of categories in sub-action class  $sa$  (as in Table 2). Negative log likelihood is used as the loss function for the captioning task

$$\mathcal{L}_{Cap} = -\frac{1}{N} \sum_{i=1}^N \sum_{sl} \ln(x_{i,sl}^{cap}) \quad (3)$$

with  $sl$  is the sentence length. The overall objective function to be minimized is the summation of all the losses

$$\mathcal{L}_{MTL} = \alpha \mathcal{L}_{AQA} + \beta \mathcal{L}_{AR} + \gamma \mathcal{L}_{Cap}. \quad (4)$$

where  $\alpha, \beta, \gamma$  are loss the weights. Now, we will introduce two different architectures for MTL-AQA.

**MTL-AQA architectures** Unlike action recognition that may be accomplished by looking at as little evidence as just a single frame [11], for AQA the complete action sequence needs to be considered because the athlete can make or lose points at any point during the whole sequence.

While spatio-temporal representations learnt using 3D CNN’s capture appearance and salient motion patterns [24], which makes them one of the best candidates for action recognition [24, 8] and also for AQA [16, 26, 13], 3D CNN’s require large memories which limits their application to small clips. We tackle this bottleneck in two ways:

1. divide the video (96 frames) into small clips (16 frames), and then aggregate clip-level representations to obtain video-level description (Sec. 4.1)
2. downsample the video into a small clip (Sec. 4.2)

Networks designed for multitask learning generally two segments: **common network backbone** and **task-specific heads**. Common network backbone learns shared representations, which are then further processed through task-specific heads to obtain more task-oriented features and outputs.

### 4.1. Averaging as aggregation (C3D-AVG)

The first network we present is C3D-AVG (Fig. 2).

**Network backbone:** Backbone consists of C3D network [24] upto the fifth pooling layer.

**Aggregation scheme:** An athlete gathering (or losing) points throughout the action can be seen as an addition operation. Combining this perspective with a good rule of thumb that when good representations are learned, linear operations on them become meaningful, we propose to enforce a linear combination of representations to be meaningful, in order to learn good representations. Specifically,



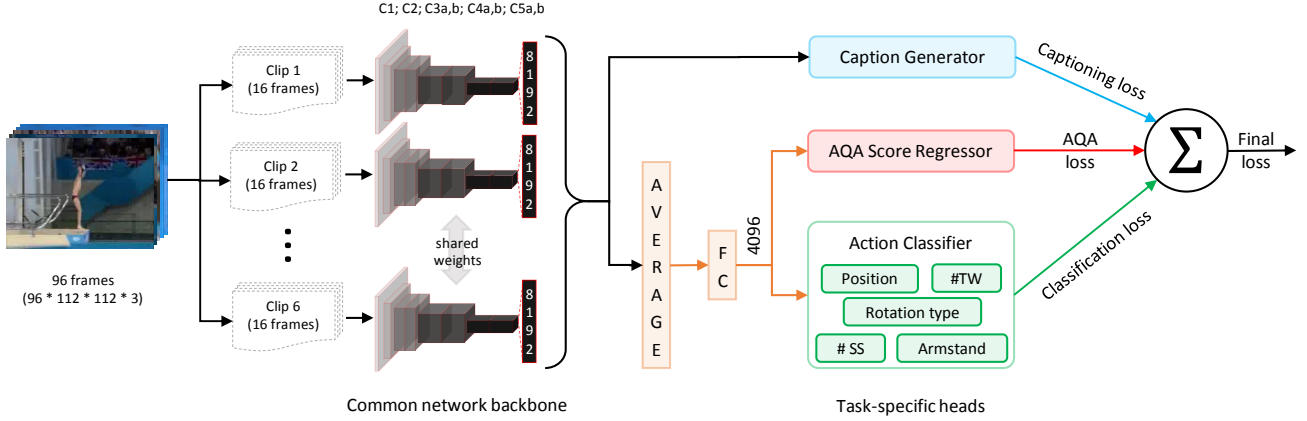


Figure 2: C3D-AVG-MTL network.

we propose to use *averaging* as the linear combination. The network is optimized end-to-end for all three tasks.

C3D-AVG network up to *Average* layer can be considered as an encoder, which encodes input video-clips into representations that when averaged (in feature space) would correspond to the total AQA points gathered by the athlete. Subsequent layers can be thought of decoders for individual tasks.

**Task-specific heads:** For action-recognition and AQA tasks, clip-level `pool-5` features are averaged element-wise to yield a video-level representation. Since captioning is a sequence-to-sequence task, the individual clip-level features are input to the captioning branch before averaging (individual clip-level features worked better in practice than averaged clip-level features for captioning).

#### 4.2. Multiscale Context Aggregation with Dilated Convolutions (MSCADC)

Multiscale context aggregation with dilated convolutions (MSCADC) [28] has been shown to improve the classification of dives in the work of Nibali *et al.* [15]. Given its strong performance on an auxiliary task MSCADC was selected for MTL. Our MTL variant network has a backbone and multiple heads as illustrated in Table 3.

**Network backbone:** The MSCADC network is based on C3D network [24] and incorporates improvements like using Batch Normalization [9] to provide better regularization which is needed in AQA where data is quite limited. Additionally, pooling is removed from the last two convolutional groups of C3D and instead a dilation rate of 2 is used. This backbone structure is shared among all the MTL tasks.

**Task-specific heads:** We use separate heads, one for each task. Heads consist of a context net followed by a few additional layers. The context net is where the feature maps are aggregated at multiple scales.

Dilated convolutions and multi-scale aggregation have shown improvements in the tasks involving dense predictions [28]. We believe that removing pooling layers and using dilated convolutions better maintains the structure of the diving athlete without losing resolution. This helps in better assessment of the athlete’s pose which is critical for AQA. For example, pose can identify when legs are aligned or split which is useful not only for diving but also other sports such as gymnastic vault, figure skating, skiing, snowboarding, etc.

Unlike the C3D-AVG network, we downsample the complete action into a short sequence of only 16 frames (something like key action snapshots) as done by Nibali *et al.* [15]. This reduces our 96-frames videos into key action snapshots which helps in processing the complete action sequence in a single pass. Processing an action sequence using this network can be thought of as distilling information from the input frames and putting it into feature maps, with different feature maps containing different kinds of pose information. A natural benefit of downsampling the sequence is that there is a significant reduction in the the number of network parameters and memory which can be used instead to increase spatial resolution.

#### 5. Experiments

**Implementation:** PyTorch [17] is used to implement all the networks; common network backbones were pretrained on the UCF101 [22] action recognition dataset. The captioning module utilized a GRU [3] cell and a dropout rate of 0.2 in the encoder and decoder. Maximum caption length is set to 100 words. Full vocabulary size is 5779. The parameters  $\alpha$ ,  $\beta$ , and  $\gamma$  in Eq. 4 are set to 1, 1, and 0.01. All networks used the Adam optimizer [12] and were trained for 100 epochs with initial learning rate of 1e-4. Data augmentation is performed through center cropping with temporal augmentation and random horizontal flipping. The center

<i>(Common network body)</i>		
C3(32); BN		
MP(1,2,2)		
C3(64); BN		
MP(2,2,2)		
{C3(128); BN} x2		
MP(2,2,2)		
{C3(256); BN} x2		
{C3(d=2,256); BN} x2		
Dropout(0.5)		
<i>(Task-specific heads)</i>		
<i>(AQA Score Head)</i>	<i>(Action recognition Head)</i>	<i>(Captioning Head)</i>
C1(12)	C1(12)	C1(12)
{Cntxt net}	{Cntxt net}	{Cntxt net}
MP(2,2,2)	MP(2,2,2)	MP(2,2,2)
C3(12); BN	C3(12); BN	C3(12); BN
C3(1)	<i>(Action recognition sub-heads)</i>	Enc. GRU
AP(2,11,11)		Dec. GRU

Table 3: **MSCADC-MTL architecture**. C3(d,ch): 3D convolutions, ch-no. of channels, d-dilation rate. C1: 1x1x1 convolutions. BN: batch normalization. MP(kr): max pooling operation, kr-kernel size. Cntxt net: context net for multi-scale context aggregation. AP: average pooling across (2x11x11) volume.

crop was found to reliably capture both the athlete and other prominent visual cues such as splash. Batch-size was set to three samples. Additional architecture-specific implementation details are as follows:

**C3D-AVG:** The model is trained end-to-end with a  $112 \times 112$  center crop from the  $171 \times 128$  pixel input video. Each dive sample was temporally normalized to a length of 96 frames.

**MSCADC:** Since this architecture does not contain fully-connected layers and all videos are downsampled to 16 frames, there are fewer model parameters allowing the use of higher resolution video input. Frames are resized to  $640 \times 360$  pixels and  $180 \times 180$  center cropping is used.

**Evaluation metrics:** AQA is assessed using Spearman’s rank correlation, dive classification uses accuracy, and commentary uses captioning metrics of Bleu, Meteor, Rouge, and CIDEr.

### 5.1. Single-task vs. Multi-task approach

We carry out an experiment to compare the performance of STL against that of MTL. We have a total of 3 tasks: AQA, detailed action recognition, and commentary generation. This experiment first considered the STL approach to AQA task and then measured the effect of including auxiliary tasks. The evaluation is summarized in Table 4. We observe that MTL approaches perform better than STL approach for both the networks, which shows that our MTL

Tasks	C3D-AVG	MSCADC
AQA	89.60	84.72
+ Cls	89.62	85.76
+ Caps	88.78	85.47
+ Cls + Caps	<b>90.44</b>	<b>86.12</b>

Table 4: **STL vs. MTL across different architectures**. Cls - classification, Caps - captioning. First row shows STL results, while the remaining rows show MTL results.

	Nibali <i>et al.</i> [15]	Ours-MTL	
		MSCADC	C3D-AVG
<b>Position</b>	74.79	78.47	<b>96.32</b>
<b>Amstand</b>	98.30	97.45	<b>99.72</b>
<b>Rotation type</b>	78.75	84.70	<b>97.45</b>
<b># Somersaults</b>	77.34	76.20	<b>96.88</b>
<b># Twists</b>	79.89	82.72	<b>93.20</b>

Model	B1	B2	B3	B4	M	R	C
C3D-AVG	0.26	0.10	0.04	0.02	0.11	0.14	0.06
MSCADC	0.25	0.09	0.03	0.01	0.11	0.13	0.05

Table 5: **Performance on auxiliary tasks**.

approach is not limited to a network but is generalizable across networks. Other thing to note here is that MTL performance improves as we incorporate more tasks. Comparing both the architectures, we find that our C3D-AVG outperforms our MSCADC for both STL and MTL, while MSCADC has the advantage of being fast and lower memory requirement than C3D-AVG. For qualitative results, refer to Table 6 and supplementary material.

Performance on the auxiliary tasks is presented in Table 5. To the best of our knowledge there is only one work (by Nibali *et al.* [15]) on detailed dive classification. Our C3D-AVG-MTL performed best on the classification task as well. We also give captioning metrics for the two networks though there is no baseline for comparison in literature.

Secondly, we compare our models with the existing methods in Table 7. We obtain the results for all of the existing methods on our dataset. C3D-SVR was the best performing method in [16] but it does not seem to benefit from the increased number of training samples. In [16], C3D-LSTM was reported to be performing worse than C3D-SVR due to insufficient amount of training data and does outperform C3D-SVR with the expanded training data. Our MSCADC-STL works better than most of the existing methods, whereas our C3D-AVG-STL is better performing than all the existing methods. Furthermore, C3D-AVG-MTL with 90.44 correlation achieves new state-of-the-art results.


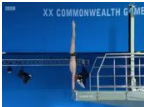
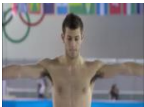

	<p><b>True labels:</b> 89.08; [Tuck, No, Backwards, 3.5, 0] <b>C3D-AVG labels:</b> 80.41; [Tuck, No, Backwards, 3.5, 0]</p> <p><b>C3D-AVG:</b> that's good she will certainly keep that with a dive that we can certainly do it in that field very impressive it is very good but it was a good dive here playing a little bit low water which is a strong start one look at that closes vertical position it's gonna get seven and a half's I think super slow-motion just throwing up a little bit too much splash but she'll get over 60s mid sixties probably and that will be good enough to keep her middle of the pack which is where she needs to be to</p>
	<p><b>True labels:</b> 63.07; [Free, Yes, Backwards, 2, 2.5] <b>C3D-AVG labels:</b> 65.79; [Free, Yes, Backwards, 2, 1.5]</p> <p><b>C3D-AVG:</b> well it's okay and the entry into the water not quite 100% vertical but he's just a little bit overcooked on the end obviously there's a few of the divers have you use themselves so if you're getting your hands out there we are rocking and rolling a little bit of a splash with technically a little bit of splash that's not the splash means that the judges will penalize him or only got to 17 from the two and a half somersaults before he goes into the water now that was a</p>
	<p><b>True labels:</b> 84.15; [Tuck, No, Backwards, 3.5, 0] <b>C3D-AVG labels:</b> 81.94; [Tuck, No, Backwards, 3.5, 0]</p> <p><b>C3D-AVG:</b> excellent excellent dive if you might he's got a lot of divers here with their hands together for him a lot of them here and take a little bit of an angle on the entry that does good through that would not quite a way over a vertical look at that perfect angle so much better judges will like that that angle so not too many</p>
	<p><b>True labels:</b> 47.77; [Pike, No, Forwards, 2.5, 1] <b>C3D-AVG labels:</b> 53.04; [Pike, No, Forwards, 2.5, 1]</p> <p><b>C3D-AVG:</b> nice nice entry because the execution was fine and then just suggesting she went surfing over the end of the diving board anyway she's a safe distance from the diving board so that's a good dive in the prelims you can see the splash moving away from the diving board six and a half's sevens at best moving further away from the podium dive after dive star with a 58 and it with a 64 this</p>

Table 6: **Qualitative results.** Labels are ordered as follows: AQA score; [Position, Armstand?, Rotation type, #SS, #TW]. Due to space constraints only generated captions are shown here; please refer to supplementary material for groundtruth.

Method	Sp. Corr.
Pose+DCT [18]	26.82
C3D-SVR [16]	77.16
C3D-LSTM [16]	84.89
Ours MSCADC-STL	84.72
Ours C3D-AVG-STL	<b>89.60</b>
Ours MSCADC-MTL	86.12
Ours C3D-AVG-MTL	<b>90.44</b>
<i>Segment-specific methods (train/test on UNLV Dive [16])</i>	
S3D (best performing in [26])	86.00
Li <i>et al.</i> [13]	80.09
Ours MSCADC-STL	79.79
Ours C3D-AVG-STL	83.83
Ours MSCADC-MTL	80.60
Ours C3D-AVG-MTL	<b>88.08</b>

Table 7: **Performance comparison with the existing AQA approaches.**

Method proposed by Xiang *et al.* [26] requires manual annotation to mark end points of all the segments which is not available in the new Diving-MTL data. Xiang *et al.* [26] used the UNLV-Dive dataset [16] so for a fair comparison with [26] we train and test our models on UNLV-Dive [16]. The results are enumerated in Table 7. Our C3D-AVG-STL does not perform as well S3D [26]. However, our C3D-AVG-MTL outperforms the S3D model. An important thing to note here is that UNLV-Dive dataset is quite a bit smaller than our newly introduced MTL-AQA dataset which should

# samples	1059	450	280	140
STL	89.60	77.27	69.63	64.17
MTL	<b>90.44</b>	<b>83.52</b>	<b>72.09</b>	<b>68.16</b>

Table 8: **STL vs. MTL generalization.** Training using increasingly reduced no. of training samples.

limit MTL performance. However, as pointed out in Section 4, MTL provides better generalization than STL, which allows C3D-AVG-MTL to learn effectively from fewer training samples.

**Generalization provided by MTL:** To ascertain that MTL is providing more generalization, we train our C3D-AVG-STL and C3D-AVG-MTL models using fewer number of datapoints. Train set size and the corresponding STL/MTL performances are detailed in Table 8. We see that MTL consistently outperforms STL, and also the gap seems to widen with fewer training samples.

## 5.2. AQA-orientedness of the learned representations

We trained our networks end-to-end to learn AQA-specific feature representation rather than relying on pre-trained action-recognition oriented features (as done in [16]). However, we question if there is a utility in learning AQA-specific feature representation or are action-recognition oriented features equally good? To answer this, we follow an evaluation scheme similar to Zhang *et al.* [30], where we train linear regressors on top of all

	c1	c2	c3	c4	c5
Baseline-1	71.01	71.39	73.13	76.34	73.69
Baseline-2	72.43	70.15	70.35	57.20	37.63
C3D-AVG-MTL	<b>74.26</b>	<b>77.95</b>	<b>82.78</b>	<b>86.18</b>	<b>85.75</b>

Table 9: **Performance of fitting linear regressors on the activations of all the convolutional layers.**

	c1	c2	c3	c4	c5
<i>Train/Test events overlapping</i>					
Baseline-1	<b>41.10</b>	32.06	36.53	46.86	<b>44.78</b>
Baseline-2	37.76	42.02	37.98	44.28	38.56
C3D-AVG-MTL	38.32	<b>42.68</b>	<b>45.53</b>	<b>49.18</b>	38.47
<i>Train/Test events non-overlapping (requires more generalization)</i>					
Baseline-1	-02.68	00.75	-03.91	-02.22	03.17
Baseline-2	-07.52	-02.44	05.07	24.09	<b>25.80</b>
C3D-AVG-MTL	-07.75	-02.77	<b>23.51</b>	<b>29.56</b>	-03.25

Table 10: **Performance of fitting linear regressors on the activations of all the convolutional layers for a novel action class, Gymnastic vault.** Top rows: Within-dataset evaluation, bottom rows: Out-of-dataset evaluation.

the convolutional layers, and compare the performance obtained for AQA and action-recognition models. In particular, we consider two action-recognition baselines: C3D model trained on UCF-101 dataset [22] (Baseline-1), and our model trained on our MTL-AQA dataset, but for factorized action recognition task (Baseline-2).

In the primary evaluation, we compare the representations for measuring the quality of diving action. Comparison is detailed in Table 9. In comparison to both the baselines, we find that our C3D-AVG-MTL learns better representations at all the intermediate layers.

Further we compare the representations for measuring the quality of an unseen action class – Gymnastic vault [16]. This helps in estimating the generalizability of the representations. We hypothesize that if our AQA network has learned better representations that actually capture the concept of *quality* in an action, then it should be able to measure the quality of an unseen action better than action-recognition specific networks. We carry out 2 different evaluations: 1) **Within-dataset evaluation** and 2) **Out-of-dataset evaluation**. In Within-dataset evaluation we randomly divide the samples into train set and test set, whereas in Out-of-dataset evaluation, train and test samples are drawn from different athletic competitions. Out-of-dataset evaluation is more challenging and requires feature representations to be more generalizable and not suffer from dataset-bias. Like the previous experiment, to com-

pare learned representations, we train linear regressors on top of all the convolutional layers. Train and test sets consist of 125 and 56 samples respectively. Results from both evaluations are presented in Table 10.

In the Within-dataset evaluation, the representations learned by all the models seem to be working well, although C3D-AVG-MTL performs best. The difference in performance becomes clearer in the Out-of-dataset evaluation. As expected, Out-of-dataset evaluation is more challenging and performances of all the models drop. However, the performances of Baseline-2 and our model drop more gracefully.

## 6. Discussion

We introduced a multitask learning approach to AQA and showed that MTL performs better than STL because of better generalization which is especially important in AQA and skill assessment since datasets are small. We showed that the representations learned by our MTL models are better able to capture the inherent concept of quality of actions. Our approach is scalable since the supervision required for the auxiliary tasks is readily available from the existing video footage with minimal extra effort compared to just AQA labeling. In addition, state-of-the-art performance was achieved without any finetuning of hyperparameters. Our best performing and recommended model, C3D-AVG-MTL, achieved 90.44% correlation with judged scores which still leaves a small gap to achieve human-experts-level performance (96% [18]).

**Extension to other actions and skills assessment:** Although this paper is geared specifically toward multitask diving AQA, the approach is general in nature. No design decisions were biased towards or specific to the diving tasks. Experiments even showed that the models trained on diving do work reasonably well for another action, gymnastic vault. This encouraging result hints at the direct application of our MTL approach on other actions and everyday skills assessment. Commentary and action class details are available almost all the of time in the sport footages. For non-sport skills assessment, such as surgery, needle passing, drawing, or painting, experts could be used to generate comments and definition of sub-actions for classification. Note that existing datasets can simply be augmented to include additional labels, instead of building new datasets from scratch. Also, our MTL approach is complementary to the existing AQA and skills assessment approaches.

**Acknowledgements:** Thank you Andy (Squadra), Mark (Wilbourne), Josh (Rana) for helping us with the dataset collection!



## References

- [1] Gedas Bertasius, Hyun Soo Park, X Yu Stella, and Jianbo Shi. Am i a baller? basketball performance assessment from first-person videos. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2196–2204. IEEE, 2017. [2](#)
- [2] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997. [4](#)
- [3] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014. [5](#)
- [4] Hazel Doughty, Dima Damen, and Walterio Mayol-Cuevas. Who’s better? who’s best? pairwise deep ranking for skill determination. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [1](#), [2](#)
- [5] Hazel Doughty, Walterio Mayol-Cuevas, and Dima Damen. The pros and cons: Rank-aware temporal attention for skill determination in long videos. *arXiv preprint arXiv:1812.05538*, 2018. [2](#)
- [6] Sonal Gupta, Joohyun Kim, Kristen Grauman, and Raymond Mooney. Watch, listen & learn: Co-training on captioned images and videos. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 457–472. Springer, 2008. [3](#)
- [7] Sonal Gupta and Raymond Mooney. Using closed captions to train activity recognizers that improve video retrieval. In *Proceedings of the CVPR-09 Workshop on Visual and Contextual Learning from Annotated Images and Videos (VCL)*, Miami, FL, June 2009. [3](#)
- [8] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6546–6555, 2018. [4](#)
- [9] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015. [5](#)
- [10] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Evaluating surgical skills from kinematic data using convolutional neural networks. In *International Conference On Medical Image Computing and Computer Assisted Intervention (MICCAI)*, 2018. [1](#)
- [11] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014. [4](#)
- [12] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [5](#)
- [13] Yongjun Li, Xiujuan Chai, and Xilin Chen. End-to-end learning for action quality assessment. In *Pacific Rim Conference on Multimedia*, pages 125–134. Springer, 2018. [1](#), [2](#), [4](#), [7](#)
- [14] Zhenqiang Li, Yifei Huang, Minjie Cai, and Yoichi Sato. Manipulation-skill assessment from videos with spatial attention network. *arXiv preprint arXiv:1901.02579*, 2019. [2](#)
- [15] Aiden Nibali, Zhen He, Stuart Morgan, and Daniel Greenwood. Extraction and classification of diving clips from continuous video footage. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2017 IEEE Conference on*, pages 94–104. IEEE, 2017. [3](#), [5](#), [6](#)
- [16] Paritosh Parmar and Brendan Tran Morris. Learning to score olympic events. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2017 IEEE Conference on*, pages 76–84. IEEE, 2017. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [8](#)
- [17] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. [5](#)
- [18] Hamed Pirsiavash, Carl Vondrick, and Antonio Torralba. Assessing the quality of actions. In *European Conference on Computer Vision*, pages 556–571. Springer, 2014. [1](#), [2](#), [3](#), [7](#), [8](#)
- [19] Ariadna Quattoni, Michael Collins, and Trevor Darrell. Learning visual representations using images with captions. In *Computer Vision and Pattern Recognition, 2007. CVPR’07. IEEE Conference on*, pages 1–8. IEEE, 2007. [3](#)
- [20] Ashish Sharma, Jatin Arora, Pritam Khan, Sidhartha Satapathy, Sumit Agarwal, Satadal Sengupta, Sankarshan Mridha, and Niloy Ganguly. Commbbox: Utilizing sensors for real-time cricket shot identification and commentary generation. In *Communication Systems and Networks (COM-SNETS), 2017 9th International Conference on*, pages 427–428. IEEE, 2017. [3](#)
- [21] Rahul Anand Sharma, K Pramod Sankar, and CV Jawahar. Fine-grain annotation of cricket videos. In *Pattern Recognition (ACPR), 2015 3rd IAPR Asian Conference on*, pages 421–425. IEEE, 2015. [3](#)
- [22] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. [5](#), [8](#)
- [23] Mohak Kumar Sukhwani. *Understanding and Describing Tennis Videos*. PhD thesis, International Institute of Information Technology Hyderabad, 2016. [3](#)
- [24] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015. [2](#), [4](#), [5](#)
- [25] Vinay Venkataraman, Ioannis Vlachos, and Pavan K Turaga. Dynamical regularity for action analysis. [1](#), [2](#)
- [26] Xiang Xiang, Ye Tian, Austin Reiter, Gregory D Hager, and Trac D Tran. S3d: Stacking segmental p3d for action quality assessment. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 928–932. IEEE, 2018. [1](#), [2](#), [4](#), [7](#)
- [27] Chengming Xu, Yanwei Fu, Bing Zhang, Zitian Chen, Yungang Jiang, and Xiangyang Xue. Learning to score the fig-

- ure skating sports videos. *arXiv preprint arXiv:1802.02774*, 2018. [2](#)
- [28] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. [5](#)
- [29] Huanyu Yu, Shuo Cheng, Bingbing Ni, Minsi Wang, Jian Zhang, and Xiaokang Yang. Fine-grained video captioning for sports narrative. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [2](#), [3](#)
- [30] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European Conference on Computer Vision*, pages 649–666. Springer, 2016. [7](#)
- [31] Aneeq Zia and Irfan Essa. Automated surgical skill assessment in rmis training. *International journal of computer assisted radiology and surgery*, 13(5):731–739, 2018. [1](#)
- [32] Aneeq Zia, Yachna Sharma, Vinay Bettadapura, Eric L Sarin, and Irfan Essa. Video and accelerometer-based motion analysis for automated surgical skills assessment. *International journal of computer assisted radiology and surgery*, 13(3):443–455, 2018. [1](#)
- [33] Aneeq Zia, Yachna Sharma, Vinay Bettadapura, Eric L Sarin, Thomas Ploetz, Mark A Clements, and Irfan Essa. Automated video-based assessment of surgical skills for training and evaluation in medical schools. *International journal of computer assisted radiology and surgery*, 11(9):1623–1636, 2016. [1](#), [2](#)