

Hand Hygiene Assessment via Joint Step Segmentation and Key Action Scorer

Chenglong Li, Qiwen Zhu, Tubiao Liu, Jin Tang, and Yu Su

Abstract—Hand hygiene is a standard six-step hand-washing action proposed by the World Health Organization (WHO). However, there is no good way to supervise medical staff to do hand hygiene, which brings the potential risk of disease spread. Existing action assessment works usually make an overall quality prediction on an entire video. However, the internal structures of hand hygiene action are important in hand hygiene assessment. Therefore, we propose a novel fine-grained learning framework to perform step segmentation and key action scorer in a joint manner for accurate hand hygiene assessment. Existing temporal segmentation methods usually employ multi-stage convolutional network to improve the segmentation robustness, but easily lead to over-segmentation due to the lack of the long-range dependence. To address this issue, we design a multi-stage convolution-transformer network for step segmentation. Based on the observation that each hand-washing step involves several key actions which determine the hand-washing quality, we design a set of key action scorers to evaluate the quality of key actions in each step. In addition, there lacks a unified dataset in hand hygiene assessment. Therefore, under the supervision of medical staff, we contribute a video dataset that contains 300 video sequences with fine-grained annotations. Extensive experiments on the dataset suggest that our method well assesses hand hygiene videos and achieves outstanding performance.

Index Terms—Hand Hygiene, Action Assessment, Step Segmentation, Key Action, Joint Learning.

I. INTRODUCTION

IN 2005, the World Health Organization designated October 15-th as “World Hand-washing Day”. But not many people can wash their hands well in life, and quite a few people have not developed good hand-washing habits. Taking novel coronavirus as an example, since the outbreak, more than 500 million people have been diagnosed and more than 6 million people have died in the world. The droplets sprayed by patients with the virus will not only spread into the air but also stick to their hands and survive for a long time. If they touch other persons or other things with their hands again, they will have a

This work is partly supported by National Natural Science Foundation of China (No. 62076003), and Anhui Provincial Key Research and Development Program (No. 202104d07020008).

C. Li is with Information Materials and Intelligent Sensing Laboratory of Anhui Province, Anhui Provincial Key Laboratory of Multimodal Cognitive Computation, School of Artificial Intelligence, Anhui University, Hefei 230601, China. (Email: lc11314@foxmail.com)

Q. Zhu, T. Liu and J. Tang are with Information Materials and Intelligent Sensing Laboratory of Anhui Province, Key Lab of Intelligent Computing and Signal Processing of Ministry of Education, Anhui Provincial Key Laboratory of Multimodal Cognitive Computation, School of Computer Science and Technology, Anhui University, Hefei 230601, China. (Email: zqw_0327@126.com, 1987055525@qq.com, tangjin@ahu.edu.cn)

Y. Su is with Hefei Normal University and Institute of Artificial Intelligence, Hefei Comprehensive National Science Center, Hefei 230061, China. (Email: yusu@hfnu.edu.cn)

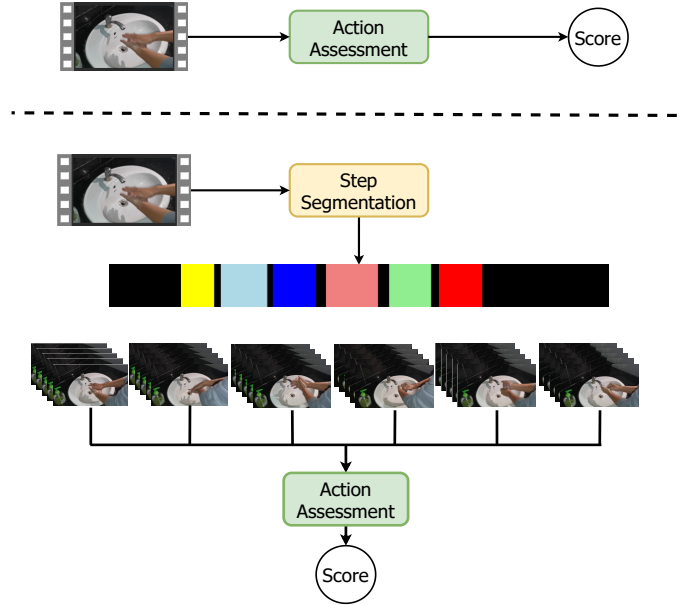


Fig. 1. Comparison of traditional action assessment framework with our framework for hand hygiene assessment. The existing framework [1], [2] directly predicts a score for a video, while we first use the step segmentation model to segment each step in the video, and then accurately assess them one by one by evaluating their key actions.

great chance of infecting other persons. Therefore, the hand is also an important medium for these viruses to spread diseases. Research shows that scientific hand washing can reduce the risk of illness by 20%. Therefore, it is very important to have a correct assessment for the hand-washing process. With the assessment, medical and other staff can correct their hand-washing actions to reduce the risk of disease spread as much as possible.

In recent years, action assessment, which aims to evaluate the quality of action, has attracted extensive attention [3]–[9]. It has important applications in many fields in the real world, such as sports and medical treatment. Previous action assessment models [1], [2] often directly evaluate a video and compute a score, as the top in Fig. 1. However, these methods would ignore many details in long actions which often involve several short steps, and thus has degraded performance in long action assessment.

The hand hygiene stipulated by the World Health Organization (WHO) is a standard long action, which includes six steps. There are three major issues in the task of hand hygiene assessment. First, each video contains some or all of the six steps, and thus existing action assessment models could not

know which frames belong to which step, which brings a big challenge to accurately assess the whole video. Second, in each step, there are several key actions to determine the quality, and thus not all actions are useful for hand hygiene assessment. Finally, there is currently a lack of a unified and high-quality annotated video dataset, which limits the research and development in hand hygiene assessment. Since there are many concepts in hand hygiene assessment, we summarize them in Table I.

TABLE I
CONCEPT EXPLANATIONS IN HAND HYGIENE ASSESSMENT.

Concept	Explanation
Long action	The whole action includes several steps and the total time is more than 10 seconds, such as hand hygiene.
Short action	The whole action includes fewer steps and the total duration is less than 10 seconds, such as diving.
Step	The certain stage in the hand hygiene assessment task.
Segment	A set of consecutive frames belonging to a step.
Key action	The action that determines the quality of the segment in each step, such as finger crossing and hand changing in the second step.

To solve the first two issues, we propose a fine-grained learning framework based on a novel approach of joint step segmentation and key action scorer for robust hand hygiene assessment. Most of existing methods [1], [2] evaluate short action such as diving by an overall prediction, which can not perform accurate evaluation of long actions since their quality depends on the completion quality of each step. Different from the overall assessment of a whole video in existing action assessment methods, we partition the input video into the step-based video segments for fine-grained scoring. Xu et al. [8] propose to segment a set of procedures for evaluation, and each two consecutive procedures does not contain the interval, i.e., the frames not belonging to any procedure. However, each two consecutive steps in hand hygiene videos usually has the interval, and the procedure segmentation is unsuitable for our task. To this end, we design a multi-stage convolution-transformer network for accurate step segmentation in hand hygiene videos. The actions in hand hygiene videos are continuous and similar, and existing multi-stage models [10], [11] easily lead to over-segmentation due to the lack of the long-range dependence. To handle this problem, we embed the linear transformer in the multi-stage convolution network to form the multi-stage convolution-transformer network, which establishes effective long-range dependence between frames without increasing too much computation.

Accurately assessing each step is the most critical issue in hand hygiene assessment. We observe that each step involves two key actions which dominate the quality of this step, as shown in Fig. 2 and Table II. Therefore, we propose a new assessment module including six key action scorers, each of which corresponds to a step and is composed of two branches with the same structure. In particular, each branch includes a fully connected layer (FC) and a learnable Sigmoid, and different branches have independent parameters to model the characteristics of different key actions.

To provide a unified platform for hand hygiene assessment,

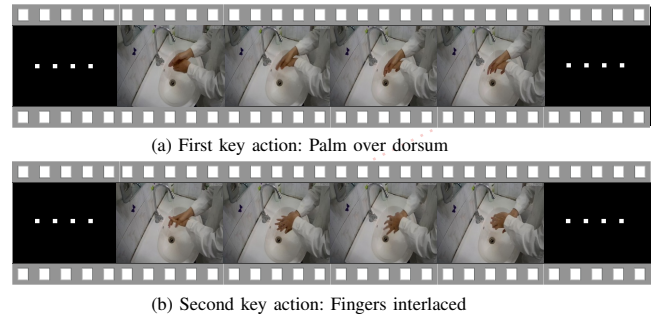


Fig. 2. Illustration of the key actions in step 2 from a hand hygiene video.

TABLE II
KEY ACTIONS IN ALL STEPS.

Step	Key action 1	Key action 2
Step 1	Palms facing each other	Finger together
Step 2	Palm over dorsum	Fingers interlaced
Step 3	Palm to palm	Fingers interlaced
Step 4	Flex fingers	Rub in the palm
Step 5	Hold your left thumb	Rotary rub
Step 6	Finger together	Rub fingertips in palms

we create a unified and high-quality video dataset, called HHA300. HHA300 contains 300 hand-washing videos of different persons from a wide range of viewpoints and background complexities, including all possible situations in real-world hand hygiene. To complete the task of step segmentation and action assessment at the same time, we provide two different forms of annotations for HHA300. In specific, We formulate a series of rules to ensure the fairness and quality of annotations, and score each hand hygiene video according to the rules. Besides, to provide high-quality fine-grained annotations, in addition to a total assessment score for each video, we also annotate the frame-level labels, including the six step labels stipulated by WHO and background label, under the supervision of medical staff.

We evaluate our framework on HHA300 dataset including step segmentation and action assessment. The results demonstrate that the joint step segmentation and key action scorer yields a notable performance gain against with other methods. In summary, our contributions are as follows.

- We propose a fine-grained learning framework based a novel approach of joint step segmentation and key action scorer for hand hygiene assessment.
- We design a multi-stage convolution-transformer network to establish the long-range dependence between frames without increasing too much computation for accurate step segmentation. It effectively alleviates the problem of over-segmentation in conventional multi-stage convolutional networks by embedding linear transformers.
- To accurately assess the quality of each step which involves several key actions, we design the key action scorer based on the multiple branch predictions, in which the Sigmoid functions are all learnable with independent parameters to model the characteristics of different key actions.

- We create the first unified and high-quality video dataset for hand hygiene assessment. It provides both video-level and frame-level annotations under the supervision of medical staff and thus would effectively promote the research and development of hand hygiene assessment.

II. RELATED WORK

A. Action Segmentation

Early action segmentation methods use sliding windows combined with non-maximum suppression [12] to focus on short-term dependence, but more windows lead to expensive calculation costs. Most recent methods use sequential convolution networks to pay attention to both short-term and long-term dependencies. MS-TCN [10] uses the temporal convolutional networks to aggregate the temporal information and uses multi-step TCN to better adjust the classification results. Li et al. [11] add the dilated convolution to solve the problem that some local information can not be extracted. BCN [13] uses the adaptive cascade network to distinguish difficult samples, which greatly improves the classification accuracy of difficult samples, and the temporal regularization method combined with action boundary information reduces the over-segmentation. The time receptive field of the model plays an important role in action segmentation. A large receptive field contributes to the long dependence between videos, while a small receptive field helps to capture local details. Gao et al. [14] propose to find a better combination of receptive fields through a global to local search scheme. Yi et al. [15] solve the problems encountered when introducing transformer into action segmentation task and design an efficient model based on transformer.

B. Action Assessment

As for hand hygiene assessment, Zhong et al. [16] apply an iterative engineering process to design a hand hygiene action detection system to improve food-handling safety, it uses the results of action classification to achieve hand hygiene assessment. Research on hand hygiene is more about pose estimation and detection tasks, which are different from hand hygiene assessment.

There are many researches on general action quality assessment. Existing methods usually regard it as a regression task, and the ultimate goal is to narrow the gap between the regression score and the ground truth score given by experts. Pirsiavash et al. [3] use the discrete cosine transform to encode joint trajectory as input feature and the support vector regression to construct the mapping from feature to final score. Based on the attention mechanism used by human beings, when evaluating videos, Li et al. [5] propose a spatial attention model based on a recurrent neural network, which considers the accumulated attention state from previous frames and advanced knowledge about the progress of ongoing tasks. Parmar et al. [4] prove that C3D [17] can effectively preserve the spatio-temporal information in the video, and help to improve the performance of evaluation tasks. Pan et al. [18] use I3D [19] as the backbone network to extract the spatio-temporal features, and establish a trainable joint

diagram and analyze their joint movements. Xu et al. [6] proposed a deep architecture that includes two LSTM to learn the different scale features of videos, and presented a figure skating sports video dataset. Recently, Zeng et al. [1] not only use the dynamic information in the video but also pay attention to the static pose information. By combining static and dynamic information in a graph-based context-aware attention module, the good video representation can be achieved. Using the idea of comparative learning, Yu et al. [2] reformulate the problem of action assessment as the regression relative score concerning another video. Xu et al. [8] construct a new fine-grained dataset and propose a procedure-aware approach to quantify quality differences between query and exemplar in a fine-grained way for action assessment. Jain et al. [9] proposed a action scoring system based deep metric learning that learns similarity between two videos. To solve the difficulty of manually label data, Zhang et al. [7] explored semi-supervised action assessment based adversarial learning and recover masked segment feature of an unlabeled video.

III. METHODOLOGY

This section will introduce the details of the proposed hand hygiene assessment approach via joint step segmentation and key action scorer, including video feature encoding, multi-stage convolution-transformer, key action scorer, and loss function. We first overview our framework for better readability.

A. Overview

The overall framework is shown in Fig. 3. Our framework for hand hygiene assessment mainly consists of two parts. The first part is the multi-stage convolution-transformer network for step segmentation, and the second part is the action assessment model based on key action scorer. First, the optical flow in hand hygiene video is extracted. The I3D feature extractor [19] takes RGB and optical flow information as inputs to compute the corresponding appearance and motion features. Then, these two features are concatenated and inputted into the multi-stage convolution-transformer network. In the network, each stage produces corresponding features, and we use the linear transformer to build the long-range dependence and thus enhance these features. After that, we obtain the temporal segments corresponding to each step and evaluate these temporal segments by predicting the quality scores respectively by the key action scorers. Finally, we sum up the score of each step to obtain the final assessment score.

B. Video Feature Encoding

Existing action segmentation networks [11], [13] usually only use RGB features for subsequent prediction. While the optical flow features [20], [21] are important for motion analysis including segmentation and assessment. On this basis, we further extract the optical flow features of video. Specifically, given an unprocessed hand hygiene video $X_{1:T} = \{x_t\}_{t=1}^T$ with T frames, we first extract the optical flow to obtain the optical flow $F_{1:T} = \{f_t\}_{t=1}^T$, and then use the pre-trained feature

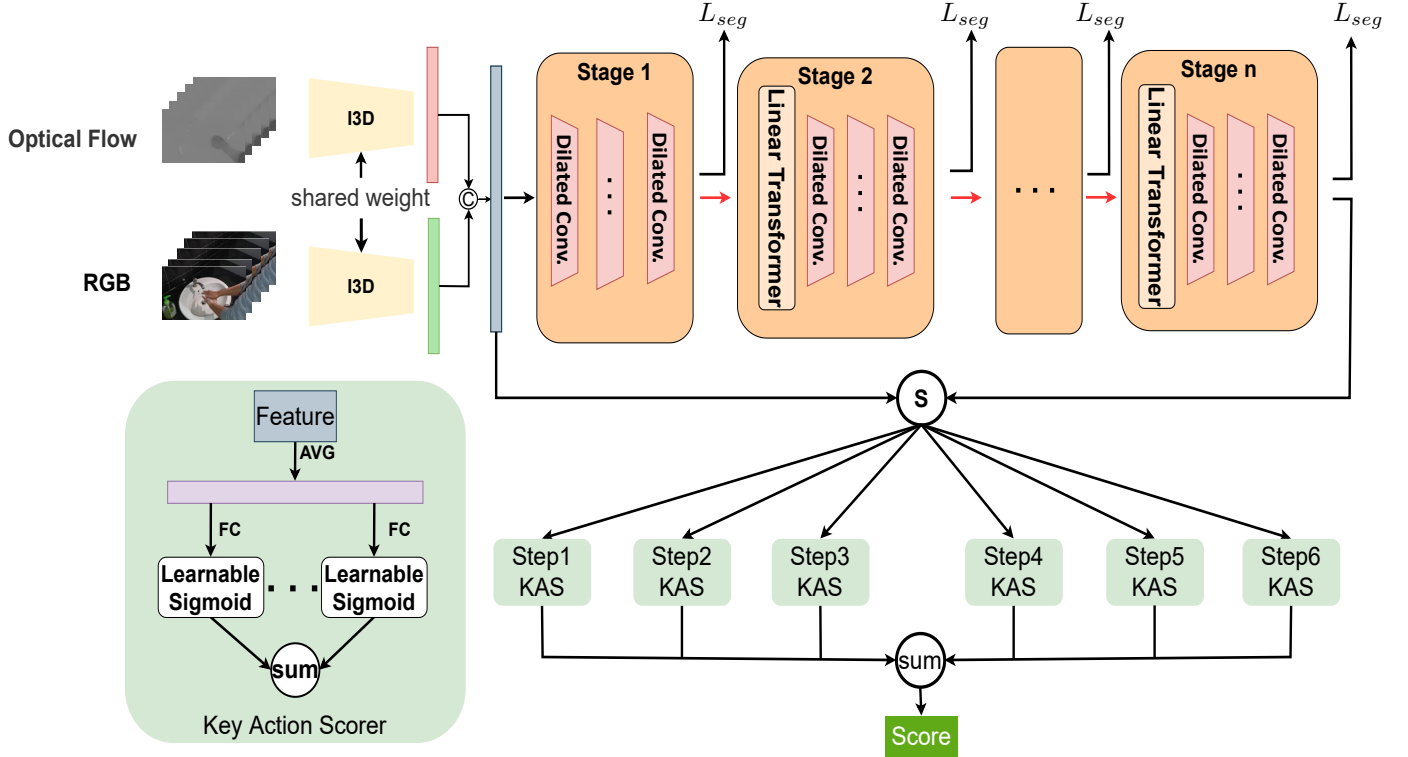


Fig. 3. Pipeline of our framework. Herein, S indicates the selection of step features, sum stands for the summation, AVG and FC are the global average pooling and the full connection layer respectively.

extraction network I3D [19] to encode RGB information and optical flow information respectively. Finally, we obtain two kinds of features $\Phi_{rgb} = \{\phi(x_1), \phi(x_1), \dots, \phi(x_T)\} \in \mathbb{R}^{T \times D}$ and $\Phi_f = \{\phi(f_1), \phi(f_1), \dots, \phi(f_T)\} \in \mathbb{R}^{T \times D}$, which represent RGB features and optical flow features respectively.

C. Step Segmentation

Most existing works [2], [8] evaluate short action such as diving, where the input is a whole video containing the short action and the output is the score of the action. Specific, given the input video feature $\Phi = \{\phi(x_1), \phi(x_1), \dots, \phi(x_T)\} \in \mathbb{R}^{T \times D}$. These works formulate action assessment task as a regression problem is to predict the action score S :

$$S = R(\Phi), \quad (1)$$

where R is the regression model. However, these methods can not perform an accurate evaluation of long action because their quality depends on the completion quality of each step. Therefore, we propose to segment the steps before evaluating the action quality. Most of existing action segmentation methods [10], [11] adopt multi-stage convolutional models. They predict a rough segmentation result in the first stage, and the result of the previous stage is gradually refined in each subsequent stage. However, the actions in the hand hygiene video are continuous and similar, which leads to segmentation errors (i.e., over-segmentation) only depending on short-range dependence in multi-stage convolutional models.

To obtain the accurate results of step segmentation, we propose to embed the linear transformer [22] in the multi-stage

convolution model [13] to form the multi-stage convolution-transformer network, which can model the long-range dependence between frames without increasing too much computation compared with the multi-stage convolutional model.

The transformer is first applied to the research of natural language processing [23]. It uses a self-attention mechanism instead of the sequence structure of RNN so that the model can be trained in parallel and get global information. The input vectors of the self-attention mechanism are usually named query, key, and value. The weight distribution of the value vector is determined by the similarity between the query vector and the key vector. Formally, the attention layer is denoted as:

$$V' = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V. \quad (2)$$

where Q , K , and V indicate the query, key, and value respectively and d_k is the vector dimension. To reduce the risk of over-segmentation caused by the lack of long-range dependence in multi-stage convolutional networks, we introduce the transformer to re-model the long-range dependence. Specifically, we design a multi-stage convolution-transformer network based step segmentation, and each stage includes several dilated convolutions with different dilated rates to establish the dependence of short-range frames. We first generate a rough result and a stage feature in the first stage using the feature extracted by I3D and then make further stage predictions. In each stage, we use the prediction result and stage feature of the previous stage to compute the refined prediction result and stage feature $\Phi_n \in \mathbb{R}^{T \times D_n}$, where n indicates the number of stages. In addition, we use the stage

feature as key, value, and query vectors at the same time, and enhance the feature of previous stage through self-attention, to establish the global dependency between all frames.

Note that the original self-attention mechanism is calculated by matrix multiplication. With the increase of sequence length T , the calculation cost is extremely expensive. Therefore, we introduce a variant of transformer, linear transformer [24], which uses another kernel function to measure correlation instead of matrix multiplication to reduce the calculation amount of transformer and has similar performance. Specifically, we use

$$\text{sim}(Q, K) = \theta(Q) \cdot \theta(K)^T. \quad (3)$$

where $\theta(x) = \text{elu}(x) + 1$. $\text{elu}()$ [25] is an improved version of ReLU function [26]. After this transformation, we only compute $\sum_{j=1}^T \theta(K_j) V_j^T$ and $\sum_{j=1}^T \theta(K_j)$ once and reuse them for every query, and the computational complexity is reduced from $O(T^2)$ to $O(T)$ as follows:

$$\begin{aligned} V'_i &= \frac{\sum_{j=1}^T \text{sim}(Q_i, K_j) V_j}{\sum_{j=1}^T \text{sim}(Q_i, K_j)} \\ &= \frac{\theta(Q_i)^T \sum_{j=1}^T \theta(K_j) V_j^T}{\theta(Q_i)^T \sum_{j=1}^T \theta(K_j)}. \end{aligned} \quad (4)$$

where i, j are the row indexes of Q and K, V respectively and $Q = W_q \Phi_n, K = W_k \Phi_n, V = W_v \Phi_n$. W_q, W_k, W_v represent the learnable linear transformation matrices.

In addition, the attention mechanism of the transformer model does not contain the positional information. To solve this problem, the transformer adds an extra vector of positional encoding to the input of the encoder layer and decoder layer. But we think that because the video itself has temporal information, the positional encoding in the transformer will disturb the temporal information, and thus we do not use positional encoding in our method. We add the linear transformer to all the later stages, and use its powerful global modeling ability to enhance the long-range dependence of the feature of each stage, so as to solve the over-segmentation problem. Finally, the step segmentation result is obtained by the last stage.

D. Hand Hygiene Assessment

After step segmentation, for each step, we choose the continuous and longest part as the representative segment. On one hand, a step might appear many times in video sequence due to possible wrong segmentation and we choose the continuous frames to filter out some wrong segmentation frames. On the other hand, after training the step segmentation model, even if there is wrong segmentation, the longest part is still correctly segmented for most frames. Therefore, we use these continuous and longest parts as the input of the module of hand hygiene assessment. After obtaining each step segment, accurately evaluating each step is the most critical issue. We observe that each step involves two key actions, which determine the quality of this step. To accurately evaluate each step in hand hygiene, we design a Key Action Scorer (KAS) to assess the key action in each step. Each step segment will input to KAS to compute the corresponding assessment score.

Key Action Scorer. For each step, there are some key actions, which play a significant role in the standardization of this step. Therefore, it is necessary to accurately assess the quality of these key actions.

Sigmoid, a common activation function, can map a variable between 0 and 1, which is widely used in different tasks. However, there are some differences among different key actions. Therefore, for each key action, a specific structure should be used to assess it. The standard Sigmoid is too steep and not suitable for the evaluation of different key actions. While we introduce a learnable parameter into Sigmoid function to control the steepness, and use learnable Sigmoid for more accurate assessment of different key actions. The Sigmoid with learnable parameter is formulated as follows:

$$LS = \frac{1}{1 + e^{-\lambda x}} \quad (5)$$

where x is the output feature of fully connected layers (FC), λ is the learnable parameter that controls the steepness of Sigmoid function.

In this work, based on the characteristics of the learnable Sigmoid function, we design a new hand hygiene assessment module including six Key Action Scorers (KAS). Each KAS corresponds to a hand hygiene step and consists of two branches with the same structure. In particular, each branch includes a FC and a learnable Sigmoid, and different branches have independent parameters to model the characteristics of different key actions. As shown in Fig. 3, we first use global average pooling to pool the segment features and then input them into different key action assessment branches which are composed of FC and learnable Sigmoid. The number of key action assessment branches is determined by the number of key actions in each step. In our work, we set the number of branches is two, as shown in Table II. We hope that each evaluation branch can learn a set of specific parameters to accurately evaluate specific key actions as follows:

$$s_i = \frac{1}{k} \sum_{j=1}^k LS_j(FC_j(\text{AVG}(\phi_i))) \quad (6)$$

where s_i, ϕ_i are the i -th step score and the segment feature respectively. k is the number of the key action. LS, FC , and AVG indicate the operations of learnable Sigmoid, fully connected layer, and global average pooling respectively. The final total score is the summation of all step scores as follows:

$$S = \sum_{i=1}^6 s_i. \quad (7)$$

where S represents the total score of the hand hygiene video.

E. Loss Function

Our loss function is composed of two parts. The first part is the loss in the step segmentation model. It includes the cross entropy loss for frame-wise classification and logarithmic probability smoothing of censored mean square error, as follows:

$$L_{CLS} = \frac{1}{T} \sum_t -\log(y_{t,c}), \quad (8)$$

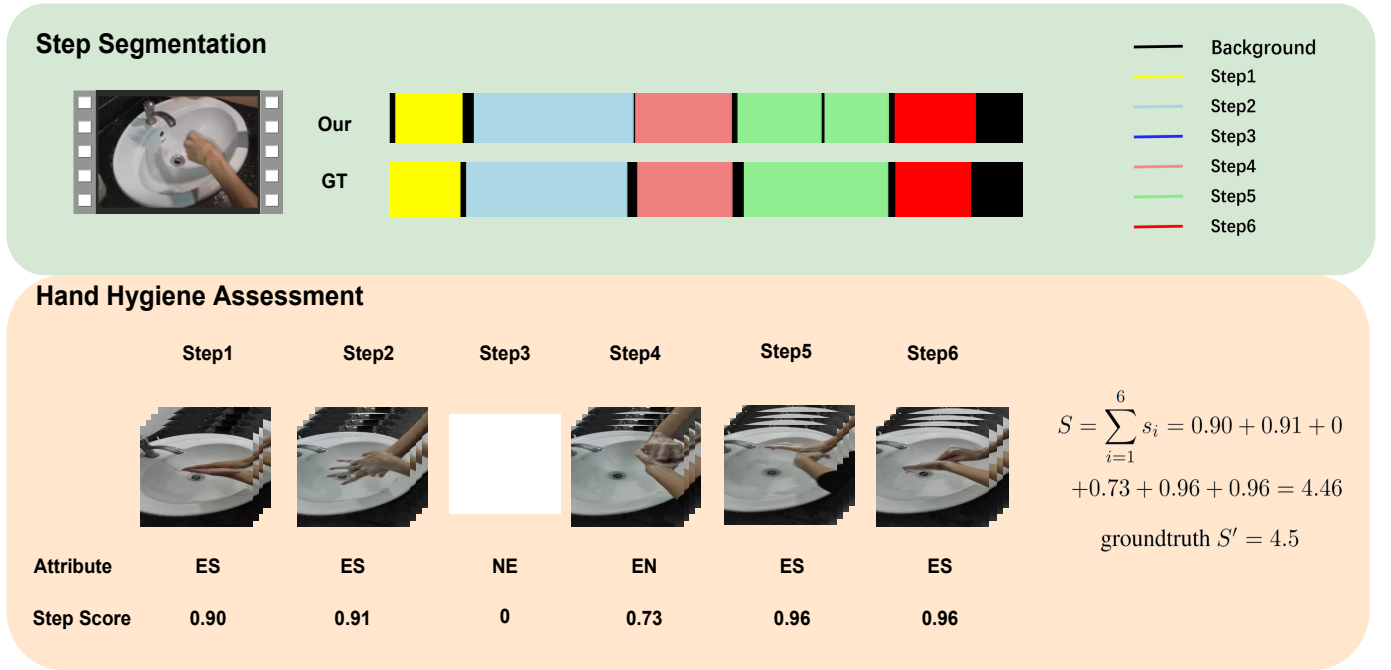


Fig. 4. An example of step segmentation and assessment by our method. The step segmentation model divides a hand hygiene video into step-based video segments. The hand hygiene assessment model assesses the key actions of each step and finally computes the final assessment score of the whole video. Herein, Attribute is the attribute annotation of each step, and Step Score indicates the output of the key action scorer.

$$\mathcal{L}_{T-MSE} = \frac{1}{TC} \sum_{t,c} \tilde{\Delta}_{t,c}^2, \quad (9)$$

$$\tilde{\Delta}_{t,c} = \begin{cases} \Delta_{t,c} & : \Delta_{t,c} \leq \tau \\ \tau & : \text{otherwise} \end{cases}, \quad (10)$$

$$\Delta_{t,c} = |\log y_{t,c} - \log y_{t-1,c}|, \quad (11)$$

where T is the video length, C is the number of classes, and $y_{t,c}$ is the probability of class c at time t . The loss in the step segmentation model as follow:

$$\mathcal{L}_{SEG} = \mathcal{L}_{CLS} + \lambda \mathcal{L}_{T-MSE}, \quad (12)$$

please refer to [10] for details.

The second part is the loss of the hand hygiene assessment, we use the mean square error to measure the difference between the predicted score and the ground truth score as follows:

$$\mathcal{L}_{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad (13)$$

where \hat{y} and y represent the predicted score and the ground truth score, respectively. The final loss function is a combination of the above mentioned losses:

$$\mathcal{L} = \mathcal{L}_{SEG} + \mathcal{L}_{MSE} \quad (14)$$

The training of the step segmentation network and the action assessment network is formulated as a multi-task learning problem and performed in an end-to-end manner. In this way, the performance of both step segmentation and action assessment is boosted. We show a typical process of our method in Fig. 4.

IV. HHA300: HAND HYGIENE ASSESSMENT DATASET

To promote the research and development of hand hygiene assessment, we create a unified dataset for hand hygiene assessment, namely HHA300. we divide HHA300 into a training set and a testing set for facilitating the evaluation of step segmentation and action assessment models, as shown in Table III. We also present some samples of HHA300 are showed in Fig. 5. In this section, we introduce the details of HHA300.

TABLE III
THE DETAILS OF OUR HHA300 DATASET.

	Video	Min frames	Mean frames	Max frames	Total frames
Train set	225	373	1048	1579	236k
Test set	75	406	1026	1436	77k

A. Data Collection

Existing hand hygiene datasets [27], [28] are all used for hand hygiene image classification. There lacks a dataset for hand hygiene assessment. To handle this problem, we create a unified hand hygiene assessment dataset HHA300, under the supervision of medical staff. In specific, we use CCD cameras to capture video data in various scenes and invite 60 persons (including ordinary people and professional medical staff) to divide into several groups to conduct hand hygiene behaviors with different standards from different scenes. Our dataset contains all most possible situations in the real world. In addition, we also invite several medical staff for professional guidance in data creation. In this way, we collect a total of 300



Fig. 5. Data samples from different viewpoints and persons in HHA300 dataset.

hand hygiene video sequences with an average video length of over 1000 frames.

B. Fine-grained Annotation

To provide high-quality fine-grained annotation, in addition to a total assessment score for each video, we also annotate the frame-level labels under the supervision of medical staff. An action segmentation and action assessment dataset needs to have high-quality frame class annotation and score annotation, which are essential for training a robust model and ensuring the fairness of performance evaluation. Therefore, we make professional annotations with the help of medical staffs, and every frame is checked to ensure the accuracy of frame-level annotation. In specific, according to the hand hygiene standard of WHO, hand hygiene is divided into six steps. We divide all frames into 7 categories, and they are palm to palm (step 1), palm over dorsum with fingers interlaced (step 2), palm to palm with fingers interlaced (step 3), back of fingers to opposing palm (step 4), rotational rubbing of the thumb (step 5), fingertips to palm (step 6) and background action. We also annotate a total assessment score for each video. To ensure high-quality annotations, we establish a series of evaluation criteria under the help of medical staffs to assess each video sequence uniformly and fairly.

C. Challenge

The complexity and type of a scene are key factors in enhancing the diversity of the dataset. To this end, We collect videos in HHA300 from a wide range of people, camera viewpoints, scene complexity, and other environmental factors. To clarify the advantages of HHA300, we analyze its diversity from the following aspects. First of all, unlike other existing datasets [27], [28], the shooting scenes are very rich, including the infirmary, public toilet, and dormitory. Secondly, we have a large number of photographers, including dozens of medical staff and students. Besides, the shooting angle is not fixed, and our dataset has six different camera viewpoints, such as top, left, and right. These factors are with different complexities and thus bring some difficulties for step segmentation and hand hygiene assessment.

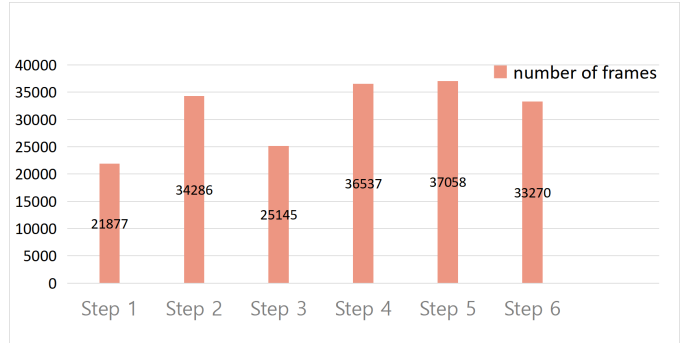


Fig. 6. The number of frames per step in HHA300 dataset.

In addition to the above challenges caused by external factors, we are also considering the challenges caused by internal factors. In the task of hand hygiene assessment, the action quality of each step has a certain influence on the result of the assessment model. Therefore, we simulate more real-world situations in data creation. In specific, for each step, we define three degrees of standardization as follows. 1) The step does not exist (NE); 2) The step exists but is nonstandard action (EN); 3) The step exists and is standard (ES). Table IV shows the video distribution of attributes on our dataset. We ask people to do hand hygiene to consider different degrees of standardization in each step. In this way, our HHA300 contains almost all possible real-world challenging situations.

TABLE IV
VIDEO DISTRIBUTION ON DIFFERENT ATTRIBUTES IN HHA300 DATASET.

Number \ Step	Step					
	1	2	3	4	5	6
Attributes						
NE	37	46	48	49	42	53
EN	27	25	32	34	21	38
ES	163	156	147	144	164	136

V. EXPERIMENTS

A. Evaluation Setting

In this section, we will describe implementation details, evaluation metrics, datasets and evaluation methods in our experiments.

Implementation. For our hand hygiene dataset HHA300, we first extract the optical flow from hand hygiene videos. Then, all frames are resized to 224×224. Finally, we extract 1024-dimensional RGB and optical flow features from the I3D model pretrained on Kinetics [19]. In all experiments, we use Adam optimizer with a learning rate of 0.0005, and the weight decay is set to 0.

Evaluation Metrics. For the evaluation of step segmentation, we employ the frame-wise accuracy (Acc), segmental edit distance, and the segmental F1 scores [29] at overlapping thresholds 10%, 25% and 50%, denoted by $F1@ \{10, 25, 50\}$. As

we all know, Acc is one of the most common evaluation metrics of step segmentation. since Acc is related to the number of frames, the influence of the long action class on this metric is greater than that of the short action class, which makes this metric unable to reflect the error of over-segmentation. To this end, two metrics are introduced, including the segmental edit distance and the segmental F1 score, which can penalize over-segmentation errors. The edit distance penalizes over-segmentation errors by predicting the order of actions, while the F1 score is similar to mean average precision (mAP) in detection task.

For the evaluation of hand hygiene assessment, similar to existing works [1], [2], we use the standard evaluation metric known as Spearman’s rank correlation coefficient ρ . which is defined as:

$$\rho = \frac{\sum_i (p_i - \bar{p})(q_i - \bar{q})}{\sqrt{\sum_i (p_i - \bar{p})^2 \sum_i (q_i - \bar{q})^2}} \quad (15)$$

where p and q represent the ranking for each sample of two series respectively. We also use the relative L2-distance(R- ℓ_2) [2] which is defined as:

$$R - \ell_2(\theta) = \frac{1}{K} \sum_{k=1}^K \left(\frac{|s_k - \hat{s}_k|}{s_{\max} - s_{\min}} \right)^2 \quad (16)$$

where s_{\min} and s_{\max} are the highest and lowest scores for an action, s_k and \hat{s} represent the ground truth score and the prediction score for the k -th sample, respectively. Spearman’s correlation focuses more on the ranks of the predicted scores while the R- ℓ_2 focuses on the numerical values of the predicted scores.

Datasets. We evaluate the proposed method on our HHA300, and further evaluate action segmentation methods on three public datasets, including **50Salads** [30], **GTEA** [31] and **Breakfast** [32] datasets. Tab V shows the comparison of these datasets.

The **50Salads** [30] dataset contains 50 videos with 17 action classes. As the name of the dataset indicates, the videos depict salad preparation activities. For evaluation, we use five-fold cross validation and report the average result.

The **Georgia Tech Egocentric Activities (GTEA)** [31] dataset contains 28 videos corresponding to 7 different activities, like preparing coffee or cheese sandwich, and is performed by 4 subjects. The frames of the videos are annotated with 11 action classes including background. On average, each video has 20 action instances. We use cross-validation for evaluation and report the average.

The **Breakfast** [32] dataset is the largest among the three datasets with 1,712 videos. The videos are recorded in 18 different kitchens showing breakfast preparation related activities. For evaluation, we use the standard 4 splits as proposed in [32] and report the average result.

Evaluation Methods. To facilitate comparison, we evaluate some step segmentation models on HHA300 dataset. In addition, following the settings for competitors given in [6], we

also evaluated some action assessment methods. we consider different combinations of the following model components.

- **Feature extractor** We use three kinds of pretrained feature extractors, namely ResNet50 [33], I3D [19], C3D [17]. For I3D and C3D, we extract RGB features and optical flow features separately. For ResNet50, we use the outputs of the fifth layer for RGB images and optical flow, and then the average pooling to obtain the features with same dimension.
- **MLP** We first use either average or maximum pooling for video-level description and then predict the score through a two-layer MLP, which is optimized by the MSE loss between the prediction and the ground truth.
- **LSTM [34]** Similar to Parmar et al. [4], we generate video-level descriptions by a LSTM architecture. In our setting, the hidden dimensions of the LSTM layers is 256, and one fully connected layer for regression.

In addition to the combination of feature extractor, MLP and LSTM, we also combine the advanced step segmentation model with MLP such as MS-TCN [10]-MLP.

B. Analysis

Hand Hygiene Assessment. We evaluate our framework against some methods on HHA300 dataset, as shown in Table VI. From the results, we can see that for the combinations of different feature extractors, MLP and LSTM, two methods based on LSTM is better than other methods based on MLP. LSTM can perform better feature descriptions by enhancing long-range dependency. But the evaluation results on all these methods are not good. The main reason is that these methods take all features of a video as the inputs of evaluation models to regress an assessment score, which does not consider fine-grained hand hygiene actions and irrelevant background actions.

For the methods based on step segmentation models and MLP, the evaluation results are better than the above methods. It is mainly because action segmentation models can divide hand hygiene into different steps for fine-grained assessment. Among these methods, BCN [13]-MLP has the best performance, mainly because BCN with boundary-aware cascade can perform more accurate step segmentation. Therefore, the performance of step segmentation is critical to hand hygiene assessment.

From the results, we can see that our framework significantly outperforms other methods in both two metrics. In particular, our method achieves 0.842 in Spearman’s rank correlation coefficient and 0.85 in relative L2-distance. Compared with other methods, our framework includes both the multi-stage convolution-transformer network for step segmentation and the key action scorer for each step of hand hygiene assessment. It not only segments long actions in hand hygiene videos into step-based segments but also makes an accurate assessment based on key actions for each step.

Step Segmentation. For the multi-stage convolution-transformer based step segmentation, we evaluate it on

TABLE V
COMPARING WITH PUBLIC ACTION SEGMENTATION DATASETS.

Datasets	Videos	Action classes	View	Description
50Salads [30]	50	17	Top-view	Salad preparation activities
GTEA [31]	28	11	Egocentric	7 Different activities, like preparing coffee or cheese sandwich
Breakfast [32]	1712	48	Third person view	Breakfast preparation related activities in 18 different kitchens
HHA300(our)	300	7	Third person view	Hand hygiene behavior of different people in different scenes

TABLE VI

RESULTS OF STEP SEGMENTATION AND HAND HYGIENE ASSESSMENT ON HHA300 DATASET. AVG STANDS FOR THE AVERAGE POOLING, MAX STANDS FOR THE MAXIMUM POOLING, MLP REFERS TO A MULTI-LAYER PERCEPTRON WITH TWO HIDDEN LAYERS, AND LSTM IS THE LONG SHORT-TERM MEMORY.

Method	F1@{10,25,50}↑	Edit↑	Acc↑	Spearman’s Correlation↑	R - ℓ_2 (*100) ↓
ResNet50-Avg-MLP	-	-	-	0.245	39.92
ResNet50-Max-MLP	-	-	-	0.281	37.22
ResNet50-LSTM	-	-	-	0.311	36.97
C3D-Avg-MLP	-	-	-	0.274	37.97
C3D-Max-MLP	-	-	-	0.286	38.54
C3D-LSTM	-	-	-	0.350	36.81
I3D-Avg-MLP	-	-	-	0.378	36.50
I3D-Max-MLP	-	-	-	0.389	36.13
I3D-LSTM	-	-	-	0.406	34.60
MS-TCN [10]-MLP	82.0 81.7 75.6	74.6	88.7	0.704	2.52
MS-TCN++ [11]-MLP	83.3 83.3 75.9	77.9	89.0	0.711	2.01
ASFR [35]-MLP	88.2 88.7 82.4	82.0	89.9	0.728	1.80
BCN [13]-MLP	87.4 87.1 81.1	81.3	89.1	0.774	1.58
Ours	89.7 89.2 83.0	83.3	89.1	0.852	1.07

HHA300 dataset shown in Table VI and compare it with four multi-stage networks. As can be seen from the results, our method is superior to other methods in three metrics. Although we are only 0.2% higher than BCN in accuracy, we are 3.7% higher in segmental edit distance and {2.3%, 3.4%, 2.7% } higher in the segmental F1 scores. The accuracy is related to the number of frames which makes this metric unable to reflect the error of over-segmentation, while the segmental edit distance and the segmental F1 scores can well reflect the ability of the model to reduce over-segmentation. Therefore, compared with other multi-stage convolution models, the experimental results show that our multi-stage convolution-transformer network can effectively alleviate over-segmentation while ensuring accuracy. Qualitative results on HHA300 dataset are shown in Fig. 7. Predictions of baseline model BCN [13] have some over-segmentation errors, but our framework can reduce these errors by the multi-stage convolution-transformer network.

In addition, we evaluate our method on three challenging datasets, including 50Salads, Georgia Tech Egocentric Activities (GTEA), and Breakfast datasets. The results are shown in Table VII. It can be seen that we achieve the state-of-the-art performance in terms of frame-wise accuracy on all datasets and competitive segmental edit distance and segmental F1 score. Compare with the baseline model BCN, we achieve 0.8%, 1.5%, and 1.1% improvements in frame-wise accuracy on 50Salads, GTEA, and Breakfast datasets respectively. We also achieve improvements in segmental edit distance and the segmental F1 score. On Breakfast dataset, our framework and BCN have similar performance in terms of segmental edit distance, but our framework outperforms it in F1@10, 25, 50 by {1.2%, 1.6%, 2.0% }. These results show that our

framework can identify the action segments that overlap with the real segments and the ground truth. However, compared with ASRF [35], although our framework is better than it in terms of frame-wise accuracy, the segmental edit distance and segmental F1 score of our framework are inferior to that of ASRF. Ishikawa et al. [35] think that correcting the prediction result by detecting the action boundary can improve segmental edit distance and segmental F1 score, but ASRF predicts action boundaries off by some margin, therefore affecting the frame-wise accuracy. Our framework can achieve a good balance in three metrics and improve other metrics without losing frame-wise accuracy.

C. Ablation Study

Effectiveness of the optical flow. Because traditional action segmentation methods [11], [13] usually only use RGB features, we introduce optical flow information to improve the discriminative ability of feature representations. Table VIII shows the comparison results of using only RGB features and RGB features with optical flow features on HHA300 dataset. As shown in the table, the introduction of optical flow information can well assist the task of step segmentation by leveraging motion cues, and thus improve all metrics.

Effectiveness of the linear transformer. We design a multi-stage convolution-transformer network based step segmentation. Table IX shows the comparison results of the multi-stage convolution network and the multi-stage convolution-transformer network on HHA300 dataset. As shown in the

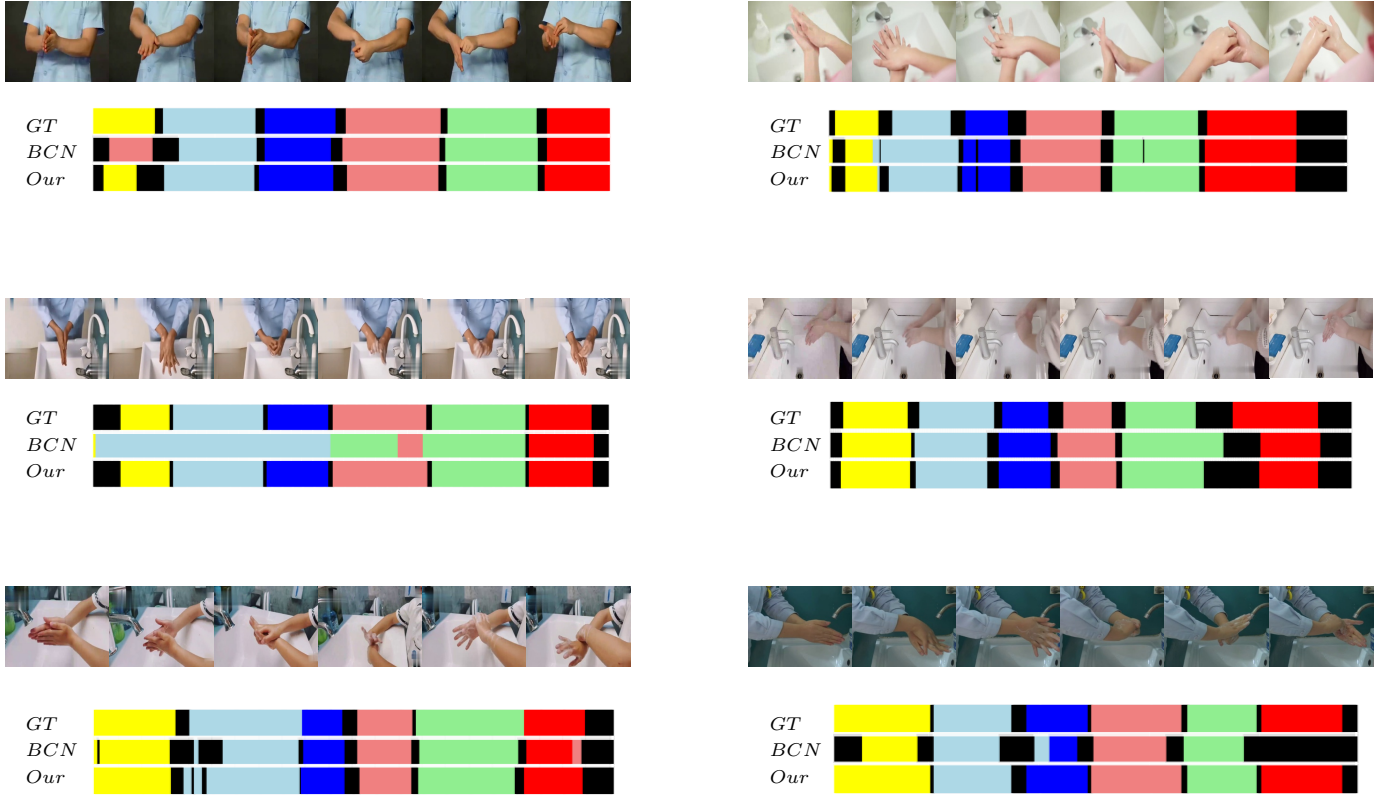


Fig. 7. Qualitative results of step segmentation on some samples from HHA300 dataset.

TABLE VII

STEP SEGMENTATION RESULTS OF OUR METHOD AGAINST OTHER METHODS ON PUBLIC DATASETS INCLUDING 50SALADS, GTEA AND BREAKFAST.

	50Salads			GTEA			Breakfast		
	Acc \uparrow	Edit \uparrow	F1@{10,25,50} \uparrow	Acc \uparrow	Edit \uparrow	F1@{10,25,50} \uparrow	Acc \uparrow	Edit \uparrow	F1@{10,25,50} \uparrow
IDT+LM [36]	48.7	45.8	44.4 38.9 27.8	-	-	-	-	-	-
ST-CNN [37]	59.4	45.9	55.9 45.6 37.1	60.6	-	58.7 54.4 41.9	-	-	-
Bi-LSTM [38]	55.7	55.6	62.6 58.3 47.0	55.5	-	66.5 59.0 43.6	-	-	-
ED-TCN [29]	64.7	59.8	68.0 63.9 52.6	64.0	-	72.2 69.3 56.0	43.3	-	-
TDRN [39]	68.1	66.0	72.9 68.5 57.2	70.1	74.1	79.2 74.4 62.7	-	-	-
SSA-GAN [40]	73.3	69.8	74.9 71.7 67.0	74.4	76.0	80.6 79.1 74.2	-	-	-
MS-TCN [10]	80.7	67.9	76.3 74.0 64.5	76.3	79.0	85.8 83.4 69.8	66.3	61.7	52.6 48.1 37.9
MS-TCN++ [11]	83.7	74.3	80.7 78.5 70.1	80.1	83.5	88.8 85.7 76.0	67.6	65.6	64.1 58.6 45.9
ASRF [35]	84.5	79.3	84.9 83.5 77.3	77.3	83.7	89.4 87.8 79.8	67.6	72.4	74.3 68.9 56.1
BCN [13]	84.4	74.3	82.3 81.3 74.0	79.8	84.4	88.5 87.1 77.3	70.4	66.2	68.7 65.5 55.0
Our	85.2	76.5	83.7 82.6 76.3	81.3	84.7	89.0 88.2 78.0	71.5	67.7	69.9 67.1 57.0

TABLE VIII

RESULTS OF OUR METHOD WITH AND WITHOUT OPTICAL FLOW FEATURES IN STEP SEGMENTATION ON HHA300 DATASET.

Method	Acc \uparrow	Edit \uparrow	F1@{10,25,50} \uparrow
Our _{rgb}	88.2	81.7	86.3 85.4 79.1
Our	89.1	83.3	89.7 89.2 83.0

TABLE IX

RESULTS OF OUR FRAMEWORK AND WITHOUT LINEAR TRANSFORMER OF STEP SEGMENTATION ON HHA300.

Method	Acc \uparrow	Edit \uparrow	F1@{10,25,50} \uparrow
Our _{w/oLT}	89.0	80.9	87.2 86.1 81.9
Our	89.1	83.3	89.7 89.2 83.0

table, after embedding the linear transformer, although the improvement in accuracy is very small, other metrics have great improvements, and these metrics are usually used to reflect the ability of the model to reduce over-segmentation. Therefore, the proposed linear transformer based multi-stage model can effectively alleviate the over-segmentation problem.

In addition, we also compare the results and computation costs of traditional transformer and linear transformer, as shown in Table X. Base on these results, we can see that the linear transformer based model can achieve similar performance with traditional transformer but has a lower computation cost.

TABLE X
RESULTS AND COMPUTATIONAL COSTS OF OUR METHOD WITH LINEAR TRANSFORMER AND WITH TRADITIONAL TRANSFORMER IN STEP SEGMENTATION ON HHA300 DATASET. $TraT$ DENOTES THE TRADITIONAL TRANSFORMER.

Method	Acc \uparrow	Edit \uparrow	F1@{10,25,50} \uparrow	params(M)	FLOPs(G)	GPU Mem.
Our $TraT$	89.8	82.7	89.4 89.1 83.1	12.46	6.14	\sim 2.63G
Our	89.1	83.3	89.7 89.2 83.0	12.46	5.77	\sim 2.09G

TABLE XI
COMPARISON OF THE STEP BASED ASSESSMENT WITH THE TRADITIONAL METHOD ON HHA300 DATASET.

Method	Sp. Corr \uparrow	R - ℓ_2 (*100) \downarrow
whole video + MLP	0.390	35.05
segment + MLP	0.852	1.07

TABLE XII
ABLATION STUDY OF KEY ACTION SCORER OF THE HAND HYGIENE ASSESSMENT ON HHA300.

Method	Sp. Corr \uparrow	R - ℓ_2 (*100) \downarrow
segment + MLP	0.814	1.49
segment + KAS(our)	0.852	1.07

Effectiveness of the step based assessment. To verify the effectiveness of our step based assessment scheme, we compare it with the traditional method that regresses the whole video to a score. To be fair, we use MLP in both methods to regress the score and the results are shown in Table XI. As can be seen from the table, our step based assessment method is far superior to the traditional method that regresses the whole video to a score. The reason is that traditional method is hard to accurately assess step based tasks, whose qualities are usually determined by key actions in each step. While we propose the step based assessment model which can solve this problem well, thus obtain much better results.

Effectiveness of the key action scorer. To verify the effectiveness of the key action scorer, we evaluate two methods in Table XII. They are the assessment method of the video after segmentation with MLP, and the assessment method of the video after segmentation with key action scorer. The results validate the effectiveness of our design.

Effectiveness of the learnable Sigmoid. We use the learnable Sigmoid instead of original Sigmoid in KAS. To verify the influence of the learnable Sigmoid on KAS, we compare it with original Sigmoid. The experimental results are shown in Table XIII. The learnable Sigmoid we selected can achieve the best performance, which demonstrates the effectiveness of the learnable Sigmoid in KAS.

VI. CONCLUSION

In this work, we present a fine-grained learning framework to perform step segmentation and key action scorer in a joint manner for accurate hand hygiene assessment. Instead of direct assessment of a whole video, we design a multi-stage

TABLE XIII
RESULTS OF OUR METHOD WITH THE LEARNABLE SIGMOID AND THE SIGMOID IN KEY ACTION SCORER (KAS) ON HHA300 DATASET.

Method	Sp. Corr \uparrow	R - ℓ_2 (*100) \downarrow
KAS(original Sigmoid)	0.821	1.42
KAS(learnable Sigmoid)	0.852	1.07

convolution-transformer network to generate step segments and then design an action assessment network based on a key action scorer to assess each step. Experimental results show that our framework can accurately assess hand hygiene videos against some state-of-the-art methods. In addition, we contribute a unified hand hygiene video dataset to promote the research and development of hand hygiene assessment. In the future, we will extend our framework to an online version for online hand hygiene assessment to improve the practicability, and expand our dataset to include more challenging scenes.

REFERENCES

- [1] L.-A. Zeng, F.-T. Hong, W.-S. Zheng, Q.-Z. Yu, W. Zeng, Y.-W. Wang, and J.-H. Lai, "Hybrid dynamic-static context-aware attention network for action assessment in long videos," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020.
- [2] X. Yu, Y. Rao, W. Zhao, J. Lu, and J. Zhou, "Group-aware contrastive regression for action quality assessment," in *Proceedings of International Conference on Computer Vision*, 2021.
- [3] H. Pirsiavash, A. Torralba, and C. M. Vondrick, "Assessing the quality of actions," in *Proceedings of European Conference on Computer Vision*, 2014.
- [4] P. Parmar and B. Tran Morris, "Learning to score olympic events," in *Proceedings of Computer Vision and Pattern Recognition Workshops*, 2017.
- [5] Z. Li, Y. Huang, M. Cai, and Y. Sato, "Manipulation-skill assessment from videos with spatial attention network," in *Proceedings of International Conference on Computer Vision Workshops*, 2019.
- [6] C. Xu, Y. Fu, B. Zhang, Z. Chen, Y.-G. Jiang, and X. Xue, "Learning to score figure skating sport videos," *IEEE transactions on circuits and systems for video technology*, vol. 30, no. 12, pp. 4578–4590, 2019.
- [7] S.-J. Zhang, J.-H. Pan, J. Gao, and W.-S. Zheng, "Semi-supervised action quality assessment with self-supervised segment feature recovery," *IEEE Transactions on Circuits and Systems for Video Technology*, 2022.
- [8] J. Xu, Y. Rao, X. Yu, G. Chen, J. Zhou, and J. Lu, "Finediving: A fine-grained dataset for procedure-aware action quality assessment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2949–2958.
- [9] H. Jain, G. Harit, and A. Sharma, "Action quality assessment using siamese network-based deep metric learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 6, pp. 2260–2273, 2020.
- [10] Y. A. Farha and J. Gall, "MS-TCN: multi-stage temporal convolutional network for action segmentation," in *Proceedings of Computer Vision and Pattern Recognition*, 2019.
- [11] S.-J. Li, Y. AbuFarha, Y. Liu, M.-M. Cheng, and J. Gall, "MS-TCN++: multi-stage temporal convolutional network for action segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [12] M. Rohrbach, S. Amin, M. Andriluka, and B. Schiele, "A database for fine grained activity detection of cooking activities," in *Proceedings of Computer Vision and Pattern Recognition*, 2012.

- [13] Z. Wang, Z. Gao, L. Wang, Z. Li, and G. Wu, "Boundary-aware cascade networks for temporal action segmentation," in *Proceedings of European Conference on Computer Vision*, 2020.
- [14] S. Gao, Q. Han, Z. Li, P. Peng, L. Wang, and M. Cheng, "Global2local: Efficient structure search for video action segmentation," in *Proceedings of Conference on Computer Vision and Pattern Recognition*, 2021.
- [15] F. Yi, H. Wen, and T. Jiang, "Asformer: Transformer for action segmentation," in *Proceedings of British Machine Vision Conference*, 2021.
- [16] C. Zhong, A. R. Reibman, H. A. Mina, and A. J. Deering, "Designing a computer-vision application: A case study for hand-hygiene assessment in an open-room environment," *Journal of Imaging*, vol. 7, no. 9, p. 170, 2021.
- [17] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *Proceedings of International Conference on Computer Vision*, 2015.
- [18] J.-H. Pan, J. Gao, and W.-S. Zheng, "Action assessment by joint relation graphs," in *Proceedings of International Conference on Computer Vision*, 2019.
- [19] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *Proceedings of Computer Vision and Pattern Recognition*, 2017.
- [20] S. Zhao, Y. Liu, Y. Han, R. Hong, Q. Hu, and Q. Tian, "Pooling the convolutional layers in deep convnets for video action recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 8, pp. 1839–1849, 2017.
- [21] E. H. P. Alwando, Y.-T. Chen, and W.-H. Fang, "Cnn-based multiple path search for action tube detection in videos," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 1, pp. 104–116, 2018.
- [22] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret, "Transformers are mns: Fast autoregressive transformers with linear attention," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5156–5165.
- [23] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [24] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret, "Transformers are mns: Fast autoregressive transformers with linear attention," in *Proceedings of International Conference on Machine Learning*, 2020.
- [25] D. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," in *Proceedings of International Conference on Learning Representations*, 2016.
- [26] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011.
- [27] H. Q. Vo, T. Do, V. C. Pham, D. Nguyen, A. T. Duong, and Q. D. Tran, "Fine-grained hand gesture recognition in multi-viewpoint hand hygiene," in *Proceedings of Systems, Man, and Cybernetics (SMC)*, 2021.
- [28] M. Ivanovs, R. Kadikis, M. Lulla, A. Rutkovskis, and A. Elsts, "Automated quality assessment of hand washing using deep learning," *arXiv preprint arXiv:2011.11383*, 2020.
- [29] C. Lea, M. D. Flynn, R. Vidal, A. Reiter, and G. D. Hager, "Temporal convolutional networks for action segmentation and detection," in *Proceedings of Conference on Computer Vision and Pattern Recognition*, 2017.
- [30] S. Stein and S. J. McKenna, "Combining embedded accelerometers with computer vision for recognizing food preparation activities," in *Proceedings of the 2013 ACM international joint conference on Pervasive and ubiquitous computing*, 2013.
- [31] A. Fathi, X. Ren, and J. M. Rehg, "Learning to recognize objects in egocentric activities," in *Proceedings of Conference on Computer Vision and Pattern Recognition*, 2011.
- [32] H. Kuehne, A. Arslan, and T. Serre, "The language of actions: Recovering the syntax and semantics of goal-directed human activities," in *Proceedings of Conference on Computer Vision and Pattern Recognition*, 2014.
- [33] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of Conference on Computer Vision and Pattern Recognition*, 2016.
- [34] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [35] Y. Ishikawa, S. Kasai, Y. Aoki, and H. Kataoka, "Alleviating over-segmentation errors by detecting action boundaries," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 2322–2331.
- [36] A. Richard and J. Gall, "Temporal action detection using a statistical language model," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3131–3140.
- [37] C. Lea, A. Reiter, R. Vidal, and G. D. Hager, "Segmental spatiotemporal cnns for fine-grained action segmentation," in *European Conference on Computer Vision*. Springer, 2016, pp. 36–52.
- [38] B. Singh, T. K. Marks, M. Jones, O. Tuzel, and M. Shao, "A multi-stream bi-directional recurrent neural network for fine-grained action detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1961–1970.
- [39] P. Lei and S. Todorovic, "Temporal deformable residual networks for action segmentation in videos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6742–6751.
- [40] H. Gammulle, S. Denman, S. Sridharan, and C. Fookes, "Fine-grained action segmentation using the semi-supervised action gan," *Pattern Recognition*, vol. 98, p. 107039, 2020.