

# Gated Forward Refinement Network for Action Segmentation

Dong Wang, Yuan Yuan\*, Qi Wang

*School of Computer Science and Center for OPTical IMagery Analysis and Learning (OPTIMAL),  
Northwestern Polytechnical University,  
Xi'an 710072, P. R. China*

---

## Abstract

Action segmentation aims at temporally locating and classifying video segments in long untrimmed videos, which is of particular interest to many applications like surveillance and robotics. While most existing methods tackle this task by predicting frame-wise probabilities and adjusting them via high-level temporal models, recent approaches classify every video frame directly with temporal convolutions. However, there are limits to generate high quality predictions due to ambiguous information in the video frames. In this paper, in order to address limitations of existing methods in temporal action segmentation task, we propose an end-to-end multi-stage architecture, Gated Forward Refinement Network (G-FRNet). In G-FRNet, each stage makes a prediction that is refined progressively by next stage. Specifically, we propose a new gated forward refinement network to adaptively correct the errors in the prediction from previous stage, where an effective gate unit is used to control the refinement process. Moreover, to efficiently optimize the proposed G-FRNet, we design an objective function that consists of a classification loss and a multi-stage sequence-level refinement loss that incorporates segmental edit score via policy gradient. Extensive evaluation on three challenging datasets (50Salads, Georgia Tech Egocentric Activities (GTEA), and the Breakfast dataset) shows our method achieves state-of-the-art results.

*Keywords:*

---

\*Corresponding author

*Email address:* y.yuan1.ieee@gmail.com (Yuan Yuan\*)

## 1 Introduction

2 Analyzing activities in videos has gained an increasing interest over re-  
3 cent years benefiting from the large amount of publicly available video data.  
4 While considerable advances have been made in human action recognition  
5 (classifying short trimmed videos) [1, 2], temporally locating and recogniz-  
6 ing actions in long untrimmed videos are still challenging. Earlier approaches  
7 [3, 4, 5] for action segmentation use sliding temporal windows of different  
8 scales to detect action segments and predict its action label. Recently, in-  
9 spired by success of neural network [6, 7, 8, 9], most methods directly predict  
10 frame-wise action labels with temporal deep models, e.g., temporal convolu-  
11 tional networks[10], or recurrent networks [11, 12]. Moreover, to capture long  
12 dependencies between video frames, the encoder-decoder architecture with  
13 stacked temporal convolutions is adapted for temporal action segmentation.  
14 Despite the advancements made by these methods, most proposed tempo-  
15 ral deep models tend to generate non-smooth predictions as they ignore the  
16 temporal continuity of action labels.

17 In order to overcome these limitations, we present a new temporal con-  
18 volutional model that stacks temporal convolutions in a multi-stage manner.  
19 Specifically, we follow the multi-stage structure in MS-TCN (multi-stage tem-  
20 poral convolutional network)[13] and take it as the baseline of our model.  
21 MS-TCN uses a sequence of dilated temporal convolutions as single-stage  
22 TCN and frame-wise features are fed into single-stage TCN to obtain an  
23 initial prediction. Moreover, to improve the accuracy of the prediction, MS-  
24 TCN sequentially stacks multiple single-stage TCNs that operate directly on  
25 the output of the previous one to get the final prediction. The effect of such  
26 composition is an implicit refinement of the predictions from the previous  
27 stages. In another words, MS-TCN simply feeds the output of the previous  
28 stage into next stage to refine the prediction. But for individual frame-wise  
29 prediction result, this implicit refinement does not tell whether it should be  
30 modified. As a result of this implicit refinement, the refined prediction from  
31 next stage may change the correct frame-wise prediction results in previous  
32 prediction and harm the subsequent refinement stage. In addition, the errors  
33 caused by implicit refinement are accumulated after several refinement stages  
34 and influence the final prediction significantly.

35 In order to alleviate above problems, we adopt the single-stage TCN from  
36 MS-TCN as the initial-stage TCN in our model and propose a new gated  
37 forward refinement network to adaptively correct the errors in previous pre-  
38 diction. Compared with implicit refinement in baseline method (MS-TCN),  
39 we introduce gate unit to control the refinement process over previous pre-  
40 diction. Specifically, the gated forward refinement network consists of one  
41 correct unit and one gate unit. Correct unit shares the same network archi-  
42 tecture with initial-stage TCN, which takes the prediction from the previous  
43 stage as input and generates the corrected results. Gate unit outputs gated  
44 weight based on hidden representations and previous prediction, which finds  
45 the errors in previous prediction by exploiting the context in the neighbor-  
46 ing labels and feature representations. And the gated forward refinement  
47 network refines the previous stage prediction in next stage according to the  
48 output of gate unit and correct unit, i.e., we refine the previous prediction  
49 with the corrected result from correct unit according to gated weight. In this  
50 way, the accumulated error from previous stages can be filtered out via small  
51 gated weight, which improves the accuracy of the refined prediction.

52 Moreover, MS-TCN applies the same loss function on every stage predic-  
53 tion, which does not explicitly require refined prediction better than previous  
54 prediction. This training objective function in MS-TCN can not force the  
55 refinement stage to correct the errors in the previous prediction. In order  
56 to tackle this problem, we introduce a multi-stage sequence-level refinement  
57 loss that forces refined prediction to achieve higher evaluation metrics than  
58 previous prediction. In detail, we integrate the traditional cross entropy loss  
59 with evaluation metric (segmental edit score) via policy gradient method,  
60 where segmental edit score difference between refined prediction and previ-  
61 ous prediction are regarded as reward in policy gradient method. In training  
62 procedure, the policy gradient method maximizes this reward to force the  
63 proposed gated forward refinement network to only correct the errors in the  
64 previous prediction. In summary, we propose a new gated forward refinement  
65 network and introduce a multi-stage sequence-level refinement loss to address  
66 the problems in MS-TCN, which results in more robust action classification  
67 and accurate action segmentation in videos.

68 We evaluate the proposed model on the following benchmark dataset-  
69 s: University of Dundee 50 Salads (50Salads) [14], Georgia Tech Egocentric  
70 Activities (GTEA) [15], and the Breakfast dataset [16]. Our results demon-  
71 strate that G-FRNet is capable of capturing dependencies between distinct  
72 actions and producing smooth predictions. Also, G-FRNet outperforms the

73 state-of-the-art on frame-wise accuracy, segmental edit score, and segmental  
74 overlap F1 score. Our key contributions include:

- 75 • We propose a new multi-stage temporal convolutional architecture with  
76 gated forward refinement network, which adaptively finds the errors in  
77 the previous prediction and corrects them based on the outputs of the  
78 proposed gate unit.
- 79 • To force the refinement network to only correct the errors in previ-  
80 ous prediction, we introduce a multi-stage sequence-level refinement  
81 loss that directly optimizes the segmental edit score via policy gradi-  
82 ent method and design a multi-stage refinement reward mechanism to  
83 suppress the errors in the predictions.
- 84 • We outperform the state of the art in action segmentation on the 50Sal-  
85 ads, GTEA and Breakfast datasets.

## 86 2. Related Work

87 Temporally locating and recognizing action segments in long untrimmed  
88 videos have been studied by many researchers. Inspired by object detection,  
89 earlier approaches [3, 4] adopt sliding window with various temporal scales  
90 to detect action segments followed with non-maximum suppression. Then,  
91 traditional methods focus on designing sophisticated models to represent  
92 actions. Fathi and Rehg [17] model actions based on the change in the  
93 state of objects and materials. In [18], an action description is constructed  
94 in terms of the interactions between hands and objects and a hierarchical  
95 model of activities, actions and objects are proposed to provide context to  
96 recognize actions. Bhattacharya et al. [19] express each video as an ordered  
97 vector time series from linear dynamical systems theory, where each time step  
98 consists of the vector formed from the concatenated confidences of the pre-  
99 trained concept detectors. These concept-based temporal representations are  
100 obtained from overlapping temporal windows and suitable for the complex  
101 event recognition. Cheng et al. [5] model temporal dependencies in the  
102 data by a sequence of visual words learnt from the video, and model the  
103 long-range dependencies by employing a Bayesian non-parametric model of  
104 discrete sequences to jointly classify and segment video sequences.

105 Other approaches follow a two-step pipeline that applies high level tem-  
106 poral modeling over frame-wise classifier. Kuehne et al. [20] combine the

107 hidden Markov model (HMM) with a context-free grammar to determine the  
108 most probable sequence of actions over frame-wise action representation. A  
109 variable-duration hidden Markov model is used in [21] to model durations of  
110 action states in addition to the transitions between action states. Vo and Bo-  
111 bick [22] address action segmentation by parsing a temporal sequence with a  
112 Bayes network, where the temporal structure of the high-level activity is rep-  
113 resented by a stochastic context-free grammar. Moreover, Richard and Gall  
114 [23] use statistical length and language modeling to represent temporal and  
115 contextual structure and find the most likely action sequence and the corre-  
116 sponding segment boundaries. While these approaches achieve good results,  
117 they are very slow as these temporal models require solving a maximization  
118 problem over very long sequences.

119 Motivated by the success of temporal convolution in speech synthesis [24],  
120 many researchers have used and modified it for the temporal action segmen-  
121 tation task. Lea et al. [10] use a hierarchy of temporal convolutions to  
122 perform fine-grained action segmentation. Their approach uses pooling and  
123 upsampling in an encoder-decoder architecture to efficiently capture long-  
124 range temporal patterns. Ding et al. [25] replace the convolutional decoder  
125 in the approach of Lea et al. [10] with a bi-directional LSTM. In addition,  
126 Lei and Todorovic [26] use deformable temporal convolutions instead of the  
127 regular temporal convolution and add a temporal residual stream for resolv-  
128 ing ambiguities about local, frame-to-frame segmentation. The approach of  
129 Farha and Gall [13] is the most related to ours, as they use a multi-stage ar-  
130 chitecture for the temporal action segmentation task. However, their model  
131 just simply stacks several identical sub-network to refine the initial predic-  
132 tion, whereas our model refines the previous predictions with a gated forward  
133 network.

134 In addition, there are a variety of different approaches that address the  
135 action segmentation task with weakly supervised action labeling [27, 28, 12,  
136 29, 30]. Kuehne et al. [29] use HMM to model the action and combine the  
137 corresponding transcripts to infer the scripted actions. They iteratively re-  
138 fine the segmentation by modifying the prediction to maximize the likelihood  
139 of all possible sequences. Following a similar pipeline, Richard et al. [30]  
140 divide each action into multiple sub-actions and propose an iterative fine-to-  
141 coarse action modeling mechanism with RNN and HMM. In the other hand,  
142 Bojanowski et al. [27] formulate the action segmentation task as a tempo-  
143 ral assignment problem and propose ordering constrained discriminative  
144 clustering to tackle it. And Huang et al. [12] propose the extended connec-

145 tionist temporal classification framework to efficiently evaluate all possible  
146 alignments via dynamic programming and explicitly enforce their consistency  
147 with frame-to-frame visual similarities. In contrast to these approaches,  
148 we address the temporal action segmentation task in a fully supervised setup  
149 and the weakly supervised case is beyond the scope of this paper.

### 150 3. Framework

151 The action segmentation is aimed at temporally locating and classifying  
152 action segments in long untrimmed videos, and we tackle it by predicting  
153 frame-wise action label for each frame. Given the frames of a video  
154  $x_{1:T} = (x_1, \dots, x_T)$ , our goal is to infer the class label for each frame  
155  $c_{1:T} = (c_1, \dots, c_T)$ , where  $T$  is the video length. The rest of this section  
156 is organized as follows. First, we give an overview of the proposed method in  
157 Section 3.1 and describe the basic-block model that stacks dilated temporal  
158 convolutions in Section 3.2, then we explain the gated forward refinement  
159 network across multiple stages in Section 3.3. Finally, we discuss the proposed  
160 multi-stage sequence-level refinement loss in Section 3.4.

#### 161 3.1. Overview

162 As shown in object detection [31], semantic segmentation [32] and other  
163 tasks, the accuracy of the initial prediction can be improved significantly by  
164 iteratively refining the prediction results. Inspired by this observation, we  
165 attempt to reduce the errors in the initial prediction by correcting them with  
166 the proposed gated forward refinement network, which iteratively refines the  
167 initial prediction over multiple steps. An overview of the proposed model is  
168 shown in Figure 1. Our model consists of multiple stages where the initial  
169 prediction from the initial stage is progressively refined with the gated forward  
170 refinement network over several steps, and the refined prediction from last  
171 step are regarded as the final prediction. Meanwhile, the proposed multi-  
172 stage sequence-level refinement loss incorporates segmental edit score via  
173 policy gradient method and a multi-stage reward mechanism is designed to  
174 guide the refinement process across the stages. Specifically, for an input  
175 video sequence  $x_{1:T}$ , each of the inputs  $x_1, \dots, x_T$  is a  $D$ -dimension feature  
176 representations extracted with pre-trained I3D network [1]. First, the feature  
177 vectors  $x_{1:T}$  are fed into the first stage network and the initial prediction  
178  $Y^0 = (y_1^0, \dots, y_T^0)$  is obtained by the initial-stage temporal convolutional  
179 network (initial-stage TCN). Then, the initial prediction  $Y^0$  is passed into the

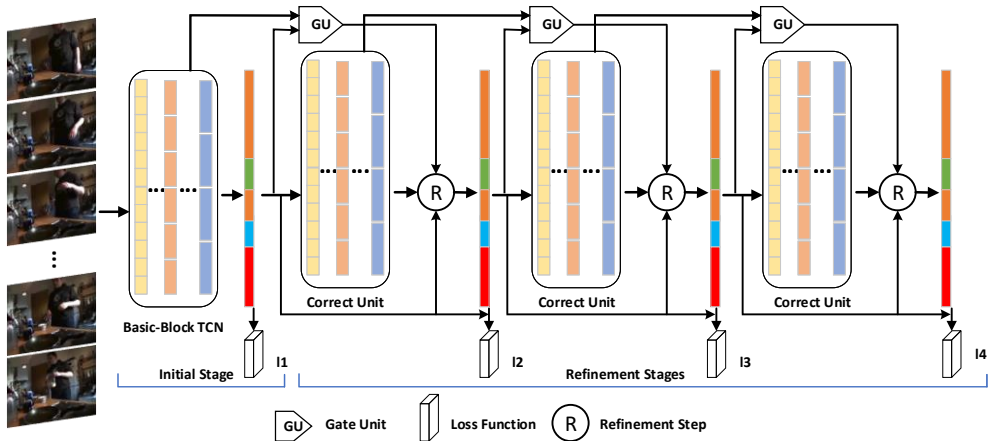


Figure 1: The overview of the proposed Gated Forward Refinement Network (G-FRNet). The frame-wise features are fed into the initial stage and obtain an initial prediction, which is progressively refined by applying several refinement stages. More details in section 3.

180 next refinement stage, which consists of basic-block temporal convolutional  
 181 network and gated forward refinement network. The refined prediction  $Y^1$   
 182 is then re-refined via the next refinement stage and all predictions after  $S -$   
 183 1 refinement stages, including the initial prediction, are denoted as  $Y^s =$   
 184  $(y_1^s, \dots, y_T^s)$ ,  $s = 0, 1, \dots, S - 1$ . Our method take the last stage prediction  
 185  $Y^{S-1}$  as the final result.

186 We propose a gated forward refinement network (G-FRNet) which form-  
 187 s the refinement stage of our model. Following previous work on multi-  
 188 stage refinement [32, 33], G-FRNet leverages feature representations from  
 189 the previous stage to obtain gated weight to correct the previous prediction.  
 190 Meanwhile, the G-FRNet generates a corrected result based on the previous  
 191 prediction and refines the previous prediction by adding the corrected re-  
 192 sult and previous prediction according to gated weight. In the other words,  
 193 by using the gated weight, the G-FRNet finds errors in previous prediction  
 194 and corrects it with the corrected result. In addition, to force the refine-  
 195 ment network to correct the errors in previous prediction, we adopt a multi-  
 196 stage sequence-level refinement loss to enforce the refined prediction from  
 197 each stage to get higher evaluation scores than the previous prediction. This  
 198 multi-stage sequence-level refinement loss integrates the segmental edit score

199 into the objective function and uses Reinforcement Learning (RL) method  
 200 to optimize it, where rewards are introduced at each stage as intermediate  
 201 supervision.

### 202 3.2. Initial-Stage TCN

203 Recently, temporal convolutional networks[34, 35] have proven their su-  
 204 perior strength in modeling the time series data[24, 36, 37]. For example,  
 205 van den Oord et al.[24] propose a fully convolutional model, called WaveNet,  
 206 for audio signal generation, and Dauphin et al.[37] introduce a convolutional  
 207 network for context dependencies modeling in language sequential analysis.  
 208 Inspired by the success of convolutional approaches in the analysis of tempo-  
 209 ral sequential data, we leverage a stack of 1D convolutional layers to model  
 210 the temporal dynamics and context dependencies over the video sequence  
 211 frames. Specifically, as shown in Figure 1, the initial stage and correct unit  
 212 in our model consist of an identical network architecture, referred as Basic-  
 213 Block TCN (temporal convolutional network), which follows the design idea  
 214 from single-stage TCN in MS-TCN.

In detail, the first layer of the initial-stage TCN is a 1D convolutional layer with kernel size 1, which takes pre-extracted features  $x_{1:T}$  or predictions  $y_{1:T}^s$  from previous stage as input. Then, the output is fed into the basic-block TCN, which stacks several the dilated residual layers and one classifier layer. Each dilated residual layer is composed of a dilated temporal convolution and a  $1 \times 1$  convolutional layer. Specifically, the output of the previous layer is first passed into dilated convolutional layer and then its output is processed by  $1 \times 1$  convolution. Moreover, residual connection is employed to facilitate gradients backpropagation. Formally, the operations at each dilated residual layer can be described as follows

$$\begin{aligned} \hat{H}_l &= ReLU(W_1 \circledast H_{l-1} + b_1) \\ H_l &= H_{l-1} + W_2 \circledast \hat{H}_l + b_2, \end{aligned} \tag{1}$$

215 where  $l \in [1, L]$  is the layer number,  $H_l$  is the output of  $l$ -th dilated residual  
 216 layer,  $\circledast$  denotes the convolution operators in dilated temporal convolution  
 217 and a  $1 \times 1$  convolutional layer.  $W_1 \in R^{3 \times D \times D}$  and  $W_2 \in R^{1 \times D \times D}$  are  
 218 learnable weights and  $b_1, b_2 \in R^D$  are bias vectors of the convolution layer.  
 219  $D$  is the number of convolutional filters.

In order to exploit the long-term temporal dependencies in the video, we stack several dilated residual layers to make the basic-block TCN cover a



large receptive field. On the top of the last dilated residual layer, to get the prediction probabilities of action class, a  $1 \times 1$  convolution followed by a SoftMax activation (classifier layer) is applied over the output of the last dilated residual layer, *i.e.*,

$$Y = \text{softmax}(W \otimes H_L + b), \quad (2)$$

220 where  $Y = y_{1:T}$  are the class probabilities for the sequence,  $H_L$  is the output  
 221 of the last dilated residual layer,  $W \in R^{C \times D}$  and  $b \in R^C$  are learnable weight  
 222 and bias of the  $1 \times 1$  convolutional layer, where  $C$  is the number of classes.

### 223 3.3. Gated Forward Refinement Network

224 Refining the initial predictions iteratively has shown significant improve-  
 225 ments in many tasks like object detection [31] and semantic segmentation  
 226 [32]. The idea of these multi-stage architectures is solving a complex prob-  
 227 lem with multiple steps such that each step only needs to make a refinement  
 228 over the output of the previous one. Another benefit of the multi-stage  
 229 method is it can effectively prevent the model from over-fitting the training  
 230 data, because it contains fewer parameters than a big model.

231 Motivated by the success of these multi-stage architectures, we adopt the  
 232 similar multi-stage structure in MS-TCN and take it as the baseline of our  
 233 model. To improve the ability of the model, in this work, we introduce a  
 234 new gate unit into refinement network to adaptively find the errors in the  
 235 previous prediction, named gated forward refinement network. Meanwhile,  
 236 a correct unit is employed to generate the corrected results from previous  
 237 prediction. The refined prediction is obtained by weighted summation of  
 238 the corrected result and previous prediction according to gate unit output.  
 239 Through the cooperation between gate unit and correct unit, the proposed  
 240 gated forward refinement network will correct the errors and keep the right  
 241 results in previous prediction.

Specifically, As shown in Figure 1, the gated forward refinement network consists of two sub-networks, one (Correct Unit) takes the prediction from the previous stage as input and generates the corrected results, another one (Gate Unit) is fed with hidden representations and prediction from the previous stage and outputs gated weight to refine the previous prediction with the corrected result from correct unit. Specifically, correct unit and gate unit contain only temporal convolutional layers, and correct unit has the identical network architecture with Basic-Block TCN. Specifically, the correct unit is

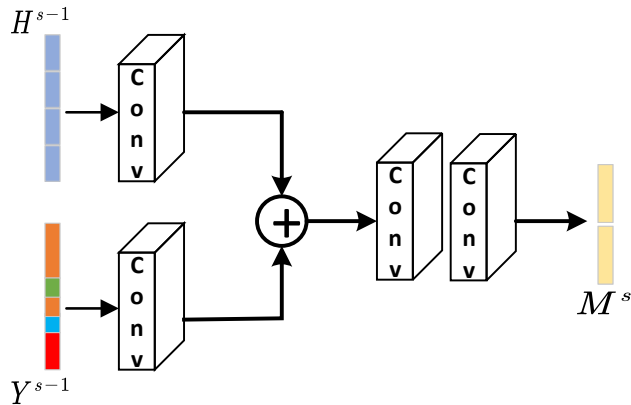


Figure 2: Overview of the gate unit.

composed of  $L$  dilated residual layers and its operations are denoted as follows

$$R^s = \mathcal{F}(Y^{s-1}), \quad (3)$$

242 where  $R^s$  is the corrected result at stage  $s$  over previous prediction,  $Y^{s-1}$  is  
 243 the output at stage  $s - 1$  and  $\mathcal{F}$  is the basic-block TCN discussed in Section  
 244 3.2. Operating on the previous prediction other than feature representa-  
 245 tion helps in capturing dependencies between action classes and generating  
 246 plausible action sequence. Furthermore, since the dimension of prediction is  
 247 far less than feature representation, there is a bottleneck layer between two  
 248 successive stages, which helps in alleviating the over-fitting problem.

**Gate Unit.** The baseline method MS-TCN [13] stacks several identical predictors sequentially to correct errors in the initial prediction progressively. Instead of generating the refined prediction with other classifier directly, we introduce gate unit to control the refinement process over the previous prediction. The gate unit is designed to find the errors in previous prediction by exploiting the context in the neighboring labels and feature representations. Figure 2 illustrates the architecture detail of the proposed gate unit. Specifically, the gate unit at stage  $s$  takes prediction  $Y^{s-1}$  and feature representation  $H_L$  (denoted as  $H^{s-1}$  for simplicity) from previous stage  $s - 1$  as its input. The features in  $H^{s-1}$  express the similarity and dissimilarity between consecutive frames, whereas class probabilities in  $Y^{s-1}$  capture the rationality of predicted action sequences. The motivation for combining  $H^{s-1}$  and  $Y^{s-1}$  in gate unit is that two consecutive timesteps with similar feature

representations should be labeled as the same action class and vice versa. A sequence of operations is carried out on  $H^{s-1}$  and  $Y^{s-1}$  followed by a softmax activation. Firstly, we apply a  $1 \times 1$  convolution with  $D$  convolutional filters to both inputs respectively, note that  $H^{s-1}$  and  $Y^{s-1}$  have different dimensions. After these operations, the two outputs are concatenated and fed into a temporal convolution with kernel size 3. Finally, another convolution layer with softmax activation is used to obtain the gated refinement weight. The formulation of operations in gate unit can be written as follows

$$\begin{aligned} g_h &= W_h \otimes H^{s-1} + b_h, \quad g_y = W_y \otimes Y^{s-1} + b_y \\ g &= W_g \otimes \text{cat}(g_h, g_y) + b_g \\ M^s &= \text{softmax}(W \otimes g + b), \end{aligned} \tag{4}$$

249 where  $W_h \in R^{1 \times D \times D}$ ,  $W_y \in R^{1 \times C \times D}$  are the convolutional weights and  $b_h$ ,  
 250  $b_y$  are the bias vectors. The outputs  $g_h, g_y \in R^{N \times T \times D}$  are concatenated  
 251 and a convolution layer with weight  $W_g \in R^{3 \times 2D \times D}$  is applied to fuse the  
 252 information from  $H^{s-1}$  and  $Y^{s-1}$ . At last, the output  $g \in R^{N \times T \times D}$  is fed into  
 253 a temporal convolution followed by softmax activation, where the learnable  
 254 weight is  $W \in R^{3 \times 2 \times D}$  and output  $M^s \in R^{N \times T \times 2}$  is gated refinement weight  
 255 to control refinement process.  $N$  is the batch size.

**Refinement Step.** Given the previous prediction  $Y^{s-1}$ , corrected result  $R^s$ , and gated refinement weight  $M^s$ , the refinement over the previous prediction can be formulated as follows

$$Y^s = M[:, :, 0] \odot R^s + M[:, :, 1] \odot Y^{s-1}, \tag{5}$$

256 where  $M[:, :, 0], M[:, :, 1] \in R^{N \times T \times 1}$  are splited from  $M$  and  $\odot$  denotes  
 257 element-wise product. Note that the sum of  $M[n, t, 0], M[n, t, 1]$  equals 1  
 258 for every  $n$  and  $t$ .

### 259 3.4. Stage-wise Supervision

The multi-stage refinement approach described above results in a deep architecture and produces a sequence of prediction. Although we are principally interested in the prediction at the last stage, predictions produced at earlier stages allow incorporating the supervision into the intermediate layers, which is validated for the vanishing gradient problem. To be specific, a two-step training procedure is designed to optimize our network. First, for prediction at each stage  $s$ , we use a combination of typical classification loss

and a regularization loss proposed in [13]. The cross entropy loss is adopted as classification loss

$$\mathcal{L}_{cls} = -\frac{1}{T} \sum_{t=1}^T \sum_{c=1}^C v_{t,c} \log(y_{t,c}), \quad (6)$$

where  $v_t \in R^C$  is the one-hot encoded vector of groundtruth action label and  $y_t \in R^C$  is the predicted probabilities vector at time  $t$ . Moreover, an additional regularization loss proposed in [13] is employed to reduce over-segmentation errors. Formally, this regularization loss is a truncated mean squared error over the frame-wise log-probabilities

$$\begin{aligned} \mathcal{L}_{T-MSE} &= \frac{1}{TC} \sum_{t,c} \tilde{\Delta}_{t,c}^2 \\ \tilde{\Delta}_{t,c} &= \begin{cases} \Delta_{t,c} & : \Delta_{t,c} < \tau \\ \tau & : otherwise \end{cases} \\ \Delta_{t,c} &= |\log y_{t,c} - \log y_{t-1,c}|, \end{aligned} \quad (7)$$

where  $T$  is the video length,  $C$  is the number of classes, and  $y_{t,c}$  is the probability of class  $c$  at time  $t$ . Note that the gradients are only computed and backpropagated through  $y_{t,c}$ . By integrating these two losses at each stage  $s$ , we obtain the first-step training objective for the full architecture by sum of the losses over all stages:

$$\mathcal{L} = \sum_s \mathcal{L}_{cls}(Y^s) + \lambda \mathcal{L}_{T-MSE}(Y^s). \quad (8)$$

260 **Multi-Stage Sequence-Level Refinement Loss.** Training with the  
 261 above loss function is prone to produce non-smooth predictions, which makes  
 262 a significant impact on action segmentation performance. The reason behind  
 263 this is that the log-likelihood score of the prediction ignores the relationship  
 264 between temporal consecutive predicted scores and does not correlate well  
 265 with the action segmentation evaluation metrics such as segmental edit score  
 266 and temporal IoU ratio. In detail, firstly, the evaluation metrics need to  
 267 identify the action segments, which is achieved by traversing the prediction  
 268 in chronological order and finding the time interval with the same action  
 269 label. Once one error appears, even surrounding with correct prediction, the  
 270 obtained action segments show serious discrepancy between the groundtruth

271 action segments. To address this problem, we incorporate the evaluation  
 272 metrics into the training objective function that punishes the non-smooth  
 273 prediction directly. Due to non-differentiability of the evaluation metrics, we  
 274 consider the action segmentation process (frame-wise action label prediction)  
 275 as a reinforcement learning problem, i.e., given an environment (the input  
 276 videos), the proposed model are viewed as an agent that conduct the frame-  
 277 wise action classification. After generating the complete prediction, the agent  
 278 will observe a sequence-level reward (evaluation metrics).

We cast our action segmenmtation model in the terminology as in [38]. Our proposed network at each stage can be viewed as an agent that interacts with external environment (input features). The policy at stage  $s$ , denoted as  $p_{\theta_s}$ , is defined by the parameters of the network  $\theta_s$ . In our model, the policy  $p_{\theta_s}$  outputs the action label sequence  $A^s$  that is obtained from the predicted class probabilities  $Y^s$  via multinomial sampling. Then, the reward is computed by an multi-stage reward mechanism that utilizes the evaluation metric (segmental edit score in this work). The goal of training is to minimize the negative expected reward over all stages:

$$\mathcal{L}_{rl}(\theta) = - \sum_s \mathcal{L}_{rl}(\theta_s) = - \sum_s \mathbb{E}_{A^s \sim p_{\theta_s}} [R(A^s)], \quad (9)$$

279 where  $A^s = \{a_0^s, \dots, a_T^s\}$  and  $a_t^s$  is sampled from the stage  $s$  at time step  $t$ .  
 280  $R(A^s)$  is the proposed multi-stage reward mechanism.

Compared with first-step training objective in Equation 8, The principal idea of our multi-stage reward mechanism is that the refined prediction at stage  $s$  should achieve higher evaluation score than the previous stage, which explicitly forces the refinement stage to correct the errors in previous prediction. Moreover, we require the result obtained at test-mode (max sampling) to be no worse than results sampled from multinomial distribution  $Y^s$ . Particularly,  $R(A^s)$  is defined as:

$$R(A^s) = \begin{cases} [r(\hat{A}^s) - r(A^s)], & : s = 0 \\ [r(\hat{A}^s) - r(A^s)] + \gamma[r(A^s) - r(A^{s-1})], & : otherwise \end{cases}, \quad (10)$$

281 where  $\hat{A}^s$  is the action label sequeunce via max sampling.  $r(\cdot)$  is the seg-  
 282 mental edit score computed by comparing the predicted action sequence to  
 283 corresponding ground-truth action sequence. In the initial stage ( $s = 0$ ),  
 284 the reward tends to supress those samples that have higher scores than max-  
 285 sampled result, because we hope the max-sampled results (same as testing)

286 achieved higher score. For the refinement stages ( $s = 1, \dots, S - 1$ ), an ad-  
 287 ditional item is used to increase the probability of the samples from stage  
 288  $s$  that outperform the samples from stage  $s - 1$  and punishes the inferior  
 289 samples.

After obtaining the reward, we use the REINFORCE algorithm [39] to compute the gradient  $\nabla_{\theta_s} L_{rl}(\theta_s)$ . REINFORCE is based on the observation that the expected gradient of a non-nondifferentiable reward function can be computed as follows:

$$\nabla_{\theta_s} L_{rl}(\theta_s) = - \sum_s \mathbb{E}_{A^s \sim p_{\theta_s}} \left[ R(A^s) \nabla_{\theta_s} \log_{p_{\theta_s}}(A^s) \right]. \quad (11)$$

In practice the expected gradient can be approximated using a single Monte-Carlo sample  $A^s = \{a_0^s, \dots, a_T^s\}$  from  $p_{\theta_s}$ , for each training example in the minibatch, the gradients are:

$$\nabla_{\theta_s} L_{rl}(\theta_s) \approx -R(A^s) \nabla_{\theta_s} \log_{p_{\theta_s}}(A^s). \quad (12)$$

290 In practice, in the second training step, we update the network by summing  
 291 the gradients from Equation 8 and Equation 12.

## 292 4. Experiments

### 293 4.1. Datasets and Experiment Setup

294 We evaluate our approach on three datasets: University of Dundee 50  
 295 Salads (50Salads) [14], Georgia Tech Egocentric Activities (GTEA) [15], and  
 296 the Breakfast dataset [16], which are widely used in action segmentation  
 297 task. The 50salads dataset contains 50 videos with 17 action classes. On  
 298 average, each video contains 20 action instances and is 6.4 minutes long.  
 299 The videos depict the salad preparation activities, which are performed by  
 300 25 actors where each actor prepares two different salads. For evaluation on  
 301 this dataset, we perform the same 5-fold cross-validation as the state-of-the-  
 302 art methods, and report the average results. The GTEA dataset contains 28  
 303 videos of seven finegrained types of daily activities in a kitchen, like preparing  
 304 coffee or cheese sandwich, performed by 4 subjects. Each video are annotated  
 305 with 11 action classes including background, showing a sequence of 20 action  
 306 instances including the background action. For this dataset, we perform the  
 307 same 4-fold cross-validation as prior work, and report the average results.  
 308 The Breakfast dataset is the largest among the three datasets with 1712

309 videos. As the name of the dataset indicates, the videos are recorded in  
310 18 different kitchens with breakfast related activities. Overall, there are 48  
311 different actions where each video contains 6 action instances on average.  
312 For evaluation, we use the standard 4 splits as proposed in [16] and report  
313 the average.

314 **Evaluation Metrics.** For all the three datasets, we use the following  
315 evaluation metrics as in [10]: frame-wise accuracy, segmental edit score and  
316 segmental overlap F1 score with threshold 10%, 25% and 50%, denoted by  
317  $F_1@{10,25,50}$ . The overlapping threshold is computed based on the inter-  
318 section over union (IoU) ratio. While the frame-wise accuracy is the most  
319 commonly used metric for action segmentation, it does not take into account  
320 the temporal structure of the prediction. For example, results with the same  
321 frame-wise accuracy may show large qualitative differences. Also, this metric  
322 does not take the temporal continuity of human actions into consideration.  
323 To address these limitations, a segmental edit score is proposed [40, 41] to  
324 penalizes this over-segmentation error. Similarly, segmental overlap F1 score  
325 also penalizes oversegmentation errors, but ignores minor temporal shifts be-  
326 tween the predictions and ground truth (which might arise from annotation  
327 noise).

#### 328 *4.2. Implementation Details and Baselines*

329 Our approach consists of one initial stage and three refinement stages  
330 that are implemented with PyTorch, i.e.,  $S = 4$ . Each basic-block TCN  
331 is composed of 10 ( $L=10$ ) dilated residual layers and the dilation factor is  
332 doubled at each layer followed by an dropout layer. We set the number of  
333 convolutional filters  $D$  to 256 in all the layer and filter size is 3. As input,  
334 following the related work [13], we first downsample the video frames and  
335 then extract I3D [1] features with 2048 dimension. For GTEA and Breakfast  
336 datasets we use the videos temporal resolution at 15 fps, while for 50Salads  
337 we downsampled the features from 30 fps to 15 fps to be consistent with  
338 the other datasets. Parameters of our network are learned using the Adam  
339 optimizer with a fixed learning rate of  $5e-5$ . The batch size  $N$  and the number  
340 of epoches are 1 and 100, respectively. We first pre-train the network with  
341 Equation 8 from scratch for 50 epoches, and the obtained network is fine-  
342 tuned by combining Equation 8 and Equation 12 for another 50 epoches.

343 **Baselines.** We take the approach of Farha and Gall [13] MS-TCN as  
344 the baseline method, which also use a multi-stage architecture for action

345 segmentation, is the most related to ours. Compared with simply stack-  
 346 ing several identical network in [13], our approach proposes an gated forward  
 347 refinement network to refine the previous prediction and an novel sequence-  
 348 level refinement loss to guide the refinement stages. Extensive evaluation  
 349 on three challenging datasets demonstrates the superiority of our approach  
 350 over [13]. For ablation studies, we specify the following G-FRNet variants:  
 351 (1) FRNet: G-FRNet without gate unit (i.e., summing  $R^s$  and  $Y^{s-1}$  directly  
 352 in Equation 5); (2) G-FRNet<sub>pre</sub>: G-FRNet trained only with Equation 8;  
 353 and (3) G-FRNet<sub>sk</sub> ( $k \in \{0, \dots, S - 2\}$ ): predicted results from the stage  
 354  $k$ . Note that the results of our approach is the prediction from the last stage  
 355  $S - 1$ . We also compare with the following closely related work: ST-CNN  
 356 [40]: Temporal convolutional filter that builds on top of spatial CNN to cap-  
 357 ture scene changes over the course of action; ED-TCN [10]: encoder-decoder  
 358 temporal convolution neural network; and TDRN [26]: two-stream network  
 359 that consists of temporal residual modules with deformable convolutions.

### 360 4.3. Experimental results

#### 361 4.3.1. Comparison with Temporal Convolution Models

362 Quantitative results on 50Salads and GTEA dataset are depicted in Table  
 363 1. We compare our G-FRNet with the most related temporal convolution  
 364 models, including ED-TCN [10], TResNet [42], TDRN [26] and MS-TCN  
 365 [13], which also is our baseline method. As shown in the table, the per-  
 366 formance of ED-TCN is worse than that of TResNet, and that TDRN is  
 367 able to improve over TResNet by computing deformable temporal convolu-  
 368 tions in two temporal streams. These methods are all single stage model  
 369 that predicts the action labels without refinement. As for the multi-stage  
 370 architecture, MS-TCN outperforms the other single stage methods on three  
 371 evaluation metrics with a large margin (up to 12.6% for frame-wise accuracy  
 372 (Acc) on the 50Salads dataset). In contrast with this baseline, our G-FRNet  
 373 achieves the significant better results in terms of segmental overlap F1 score  
 374 and segmental edit score. Specifically, our approach makes a moderate im-  
 375 provement, i.e., 2-3% in segmental overlap F1 score with different thresholds  
 376 and segmental edit score on the both datasets, demonstrating its superiority  
 377 in addressing the oversegmentation errors and non-smooth prediction result-  
 378 s. Moreover, our model achieve these improvements without increasing the  
 379 frame-wise accuracy, suggesting the advantage of correcting errors in the re-  
 380 finement stages. The effect of the refinement stage can also be seen in the  
 381 qualitative results shown in Fig 4.



Table 1: Performance comparison with respect to the most related temporal convolution models.

Method	50Salads			GTEA		
	$F_1@{\{10,25,50\}}$	Edit	Acc	$F_1@{\{10,25,50\}}$	Edit	Acc
ED-TCN [10]	68.0,63.9,52.6	59.8	64.7	72.2,69.3,56.0	-	64.0
TResNet [42]	69.2,65.0,54.4	60.5	66.0	74.1,69.9,57.6	64.4	65.8
TDRN [26]	72.9,68.5,57.2	66.0	68.1	79.2,74.4,62.7	74.1	70.1
MS-TCN [13]	76.3,74.0,64.5	67.9	80.7	85.8,83.4,69.8	79.0	76.3
<b>G-FRNet</b>	<b>78.0,76.2,67.0</b>	<b>71.4</b>	<b>80.7</b>	<b>89.1,87.5,72.8</b>	<b>83.5</b>	<b>76.7</b>

382 *4.3.2. Effect of the Refinement Stages*

383 To demonstrate the effectiveness of stacking several refinement stages  
384 over the initial stage, we compare the results from the initial stage and the  
385 results from different refinement stages, which are summarized in Table 2.  
386 As shown in the table, the initial stage TCN already achieves a comparable  
387 frame-wise accuracy, and there is no obvious improvement on frame-wise  
388 accuracy with refinement stages. Nevertheless, as we can see, the segmental  
389 edit score and segmental overlap F1 score of these predictions are of great  
390 difference, indicating the quality of these predictions are very different. To  
391 be specific, the initial prediction only gets 17.7% segmental edit score, and  
392 it increases to 48.7% after one refinement stage. On the other hand, Adding  
393 more refinement stages can progressively improve the scores and reduce  
394 the over-segmentation errors. There is a huge increase in the segmental edit score  
395 when two or three refinement stages are used, and the fourth refinement  
396 stage still improves the results but not as significantly as the previous stages.  
397 The similar improvements are shown on segmental overlap F1 score with  
398 different thresholds, indicating the refinement stage’s capacity for correcting  
399 the errors in previous prediction. Moreover, in Figure 3, we qualitatively  
400 compare the prediction from different stages on a sample test video from  
401 the 50 Salads dataset. As we can see, the errors in previous prediction are  
402 gradually corrected with the refinement stages. Also, the figure shows that  
403 the action boundaries are refined to more precise and the non-smooth results  
404 in groundtruth action segments are removed after applying more refinement  
405 stages. In another word, the sequential continuity of the action labels in one  
406 action segment and temporal dependencies between different action segments  
407 are modeled well with refinement stages. This suggests that using refinement

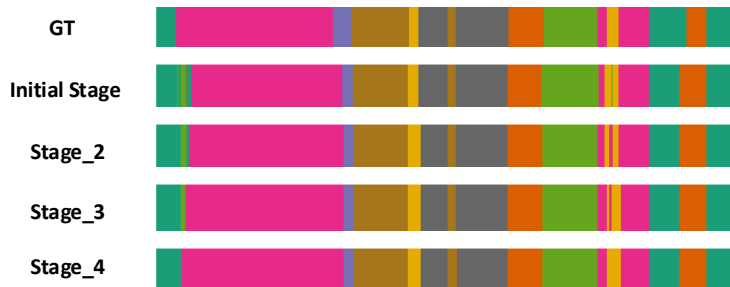


Figure 3: Qualitative result from the 50Salads dataset for comparing different predictions from different stages.

stage is critical for improving quality of prediction.

Table 2: Effect of the refinement stages on the 50Salads dataset.

	$F_1@{10,25,50}$			Edit	Acc
Initial-Stage TCN	24.1	22.2	18.6	17.9	79.3
G-FRNet (2 stages)	55.9	53.8	46.9	48.7	80.2
G-FRNet (3 stages)	70.9	69.2	60.6	63.7	80.6
G-FRNet (4 stages)	78.0	76.2	67.0	71.4	80.7

408

### 409 4.3.3. Ablation Studies

410 To verify the superiority of the proposed gated forward refinement net-  
 411 work and multi-stage sequence-level refinement loss in temporal action seg-  
 412 mentation, we compare our method with several variants baseline, including  
 413 (1) FRNet: G-FRNet without gate unit (i.e., summing  $R^s$  and  $Y^{s-1}$  directly  
 414 in Equation 5), (2) G-FRNet<sub>pre</sub>: G-FRNet trained only with Equation 8, and  
 415 MS-TCN-256 (each layer with 256 convolutional filters). All these variants  
 416 are implemented with four stages and tested on the 50Salads dataset. As  
 417 shown in Table. 4, MS-TCN-256 is worse than MS-TCN-64 in [13] on seg-  
 418 mental edit score and segmental overlap F1 score, indicating an over-fitting  
 419 problem occurs in MS-TCN method as a result of increasing the number  
 420 of parameters. Note that MS-TCN-256 achieves the best frame-wise accu-  
 421 racy, showing the benefit from increasing the parameters. In contrast, our  
 422 G-FRNet and its variants achieve comparative performance with 256 con-  
 423 volutional filters, which demonstrates the preferable generalization capacity

424 of our model.  $\text{FRNet}_{pre}$  denotes the model FRNet trained without the pro-  
 425 posed multi-stage sequence-level refinement loss.  $\text{G-FRNet}_{pre}$  outperforms  
 426 this baseline in terms of three evaluation metrics, which validates that the  
 427 proposed gate unit can adaptively control the refinement process over the  
 428 previous prediction. Additionally, FRNet and G-FRNet represent the mod-  
 429 els fine-tuned by the proposed multi-stage sequence-level refinement loss, and  
 430 the comparison FRNet *vs*  $\text{FRNet}_{pre}$  and G-FRNet *vs*  $\text{G-FRNet}_{pre}$  show the  
 431 effectiveness of the proposed refinement loss. Particularly, the segmental edit  
 432 score and segmental overlap F1 score of these models are significantly im-  
 433 proved, such as 67.0% vs 57.4% on  $F_1@{50}$ , 71.4% vs 63.1% on segmental  
 434 edit score between G-FRNet and  $\text{G-FRNet}_{pre}$ . Note that G-FRNet outper-  
 435 forms these variants baseline on all evaluation metrics, demonstrating the  
 superiority of our method in temporal action segmentation.

Table 3: Comparisons of performance with several variants of the proposed method on 50Salads dataset.

	$F_1@{10,25,50}$			Edit	Acc
MS-TCN(64)	76.3	74.0	64.5	67.9	80.7
MS-TCN(256)	55.8	63.3	48.1	51.1	81.4
$\text{FRNet}_{pre}$	68.1	66.3	57.5	60.7	79.1
FRNet	71.7	69.5	61.3	63.5	78.4
$\text{G-FRNet}_{pre}$	70.3	67.1	57.4	63.1	74.2
G-FRNet	78.0	76.2	67.0	71.4	80.7

436

#### 437 4.3.4. Comparing Different Reward Functions

438 As the core of multi-stage sequence-level refinement loss function, we de-  
 439 sign a multi-stage reward mechanism to force the refinement stage to correct  
 440 the errors in previous prediction. As shown in previous section, while this  
 441 refinement loss does not improve the frame-wise accuracy, we find this loss  
 442 produces a huge boost to the segmental edit score and segmental overlap F1  
 443 score, indicating its superiority in reducing the errors in previous prediction.  
 444 The reason behind this improvement is that the evaluation metric segmental  
 445 edit score are used to optimize the model via policy gradient, where desig-  
 446 ning reward function  $R(\cdot)$  is the most important factor. To demonstrate the  
 447 superiority of the proposed reward function in Equation 10, we compare our  
 448 method with several different reward functions that also incorporate segmen-  
 449 tal edit score into training. In detail, first, we directly takes the segmental

450 edit score of the prediction as the reward for that stage, i.e.,  $R(A^s) = r(A^s)$ ,  
 451 denoted as  $\text{Reward}_r$ . Second, we remove the first term in Equation 10 for the  
 452 stage  $s > 0$  and use the segmental edit score as the reward for initial stage,  
 453  $R(A^0) = r(A^0), R(A^s) = r(A^s) - r(A^{s-1}), s > 0$ , denoted as  $\text{Reward}_{r-r}$ .  
 454 And the proposed reward function 10 is denoted as  $\text{Reward}_{all}$ . As shown  
 455 in Table 4,  $\text{Reward}_{all}$  outperforms other methods on multiple metrics, espe-  
 456 cially segmental edit score. And the  $\text{Reward}_{r-r}$  achieve better results on all  
 457 metrics than  $\text{Reward}_r$ , which indicates that reward function comparing the  
 458 predictions from consecutive two stage can guide the model to correct the  
 459 errors. Futhermore, compared with the  $\text{G-FRNet}_{pre}$  results from Table 3, we  
 460 can find that the models trained with these three reward funtions all outper-  
 461 forms the pre-traind model, demonstrating the superiority of incorporating  
 462 the segmental edit score into training.

Table 4: Comparisons of performance with different reward function on 50Salads dataset.

	$F_1@ \{10,25,50\}$			Edit	Acc
$\text{Reward}_r$	77.0	74.6	65.7	69.7	80.0
$\text{Reward}_{r-r}$	<b>78.6</b>	76.0	66.7	70.4	80.6
$\text{Reward}_{all}$	78.0	<b>76.2</b>	<b>67.0</b>	<b>71.4</b>	<b>80.7</b>

#### 4.3.5. Comparison with the State-of-the-Art

463 In this section, we compare the proposed model to the state-of-the-art  
 464 methods on three datasets: 50Salads, Georgia Tech Egocentric Activities  
 465 (GTEA) and Breakfast datasets. The results are presented in Table 5. As  
 466 shown in the table, the proposed G-FRNet outputforms the state-of-the-art  
 467 methods on the three datasets in terms of all metrics. Specifically, our ap-  
 468 proach makes a moderate improvement on 50Salads and GTEA dataset, i.e.,  
 469 2%-4% in segmental edit score and segmental overlap F1 score, demonstratin  
 470 its superiority in reducing the non-smooth errors. As for Breakfast datasets,  
 471 our approach outperforms MS-TCN with a huge boost in items of all evalua-  
 472 tion metrics, i.e., over 10% improvement on segmental overlap F1 score with  
 473 different threshold. In addition, MS-TCN(IDT) replace the I3D features with  
 474 the improved dense trajectories (IDT) features, which are the standard used  
 475 features for the Breakfast dataset. And the improvement is not shown in  
 476 evaluation metrics, indicating the impact of the features is unstable. This  
 477

Table 5: Comparison with the state-of-the-art on 50Salads, GTEA, and the Breakfast dataset.

<b>50Salads</b>	$F_1@{10,25,50}$			Edit	Acc
IDT+LM [23]	44.4	38.9	27.8	45.8	48.7
Bi-LSTM [11]	62.6	58.3	47.0	55.6	55.7
ED-TCN [10]	68.0	63.9	52.6	59.8	64.7
TDRN [26]	72.9	68.5	57.2	66.0	68.1
MS-TCN [13]	76.3	74.0	64.5	67.9	80.7
G-FRNet	<b>78.0</b>	<b>76.2</b>	<b>67.0</b>	<b>71.4</b>	<b>80.7</b>
<b>GTEA</b>	$F_1@{10,25,50}$			Edit	Acc
Bi-LSTM [11]	66.5	59.0	43.6	-	55.5
ED-TCN [10]	72.2	69.3	56.0	-	64.0
TDRN [26]	79.2	74.4	62.7	74.1	70.1
MS-TCN [13]	85.8	83.4	69.8	79.0	76.3
G-FRNet	<b>89.1</b>	<b>87.5</b>	<b>72.8</b>	<b>83.5</b>	<b>76.7</b>
<b>Breakfast</b>	$F_1@{10,25,50}$			Edit	Acc
ED-TCN [10]	-	-	-	-	43.3
TCFPN [28]	-	-	-	-	52.0
GRU [30]	-	-	-	-	60.6
MS-TCN (I3D) [13]	52.6	48.1	37.9	61.7	66.3
MS-TCN (IDT) [13]	58.2	52.9	40.8	61.4	65.1
G-FRNet	<b>71.1</b>	<b>65.7</b>	<b>53.6</b>	<b>70.6</b>	<b>67.7</b>

478 suggests that designing suitable model is critical for improving accuracy of  
479 action segmentation. Futhermore, since our model does not use any recurrent  
480 layers, it is very fast both during training and testing.

## 481 **5. Conclusions**

482 We present a gated forward refinement network for the temporal action  
483 segmentation task. The model consists of two components: an initial-stage  
484 TCN for obtaining the initial prediction and several refinement stages for re-  
485 fining previous prediction. The refinement stage learns to correct the errors in  
486 the previous prediction with the proposed gated forward refinement network,  
487 which adaptively control the refinement process with gate unit. The experi-  
488 mental evaluations demonstrate the capability of this multi-stage architecture  
489 in capturing temporal dependencies among action classes. Moreover, in order  
490 to force the refinement network to correct the errors in previous prediction,  
491 we introduce a multi-stage sequence-level refinement loss that incorporates  
492 the non differentiable segmental edit score into training via policy gradient  
493 method, and the non-smooth and over-segmentation errors in the prediction  
494 are significantly reduced. By combining these improvements, our model out-  
495 performs the state-of-the-art methods on three challenging datasets with a  
496 large margin and extensive evaluations have demonstrated the superiority  
497 of the proposed methods in addressing non-smooth and over-segmentation  
498 errors.

## 499 **Acknowledgements**

500 This work was supported by the National Key R&D Program of Chi-  
501 na under Grant 2017YFB1002202, National Natural Science Foundation of  
502 China under Grant 61632018, 61825603, U1864204 and 61773316.

## 503 **References**

- 504 [1] J. Carreira, A. Zisserman, Quo vadis, action recognition? a new model  
505 and the kinetics dataset, in: IEEE Conference on Computer Vision and  
506 Pattern Recognition, 2017, pp. 6299–6308.
- 507 [2] C. Feichtenhofer, A. Pinz, R. Wildes, Spatiotemporal residual networks  
508 for video action recognition, in: Advances in Neural Information Pro-  
509 cessing Systems, 2016, pp. 3468–3476.

- 510 [3] M. Rohrbach, S. Amin, M. Andriluka, B. Schiele, A database for fine  
511 grained activity detection of cooking activities, in: IEEE Conference on  
512 Computer Vision and Pattern Recognition, 2012, pp. 1194–1201.
- 513 [4] S. Karaman, L. Seidenari, A. Del Bimbo, Fast saliency based pooling of  
514 fisher encoded dense trajectories, in: European Conference on Computer  
515 Vision, THUMOS Workshop, Vol. 1, 2014, p. 5.
- 516 [5] Y. Cheng, Q. Fan, S. Pankanti, A. Choudhary, Temporal sequence mod-  
517 eling for video event detection, in: IEEE Conference on Computer Vision  
518 and Pattern Recognition, 2014, pp. 2227–2234.
- 519 [6] Q. Wang, J. Gao, Y. Yuan, Embedding structured contour and location  
520 prior in siamesed fully convolutional networks for road detection, IEEE  
521 Transactions on Intelligent Transportation Systems 19 (1) (2018) 230–  
522 241.
- 523 [7] B. Zhao, X. Li, X. Lu, Hierarchical recurrent neural network for video  
524 summarization, in: ACM international conference on Multimedia, 2017,  
525 pp. 863–871.
- 526 [8] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, L. Van Gool,  
527 Temporal segment networks: Towards good practices for deep action  
528 recognition, in: European conference on computer vision, 2016, pp. 20–  
529 36.
- 530 [9] Y. Zhao, Y. Xiong, L. Wang, Z. Wu, X. Tang, D. Lin, Temporal ac-  
531 tion detection with structured segment networks, in: IEEE International  
532 Conference on Computer Vision, 2017, pp. 2914–2923.
- 533 [10] C. Lea, M. D. Flynn, R. Vidal, A. Reiter, G. D. Hager, Temporal con-  
534 volutional networks for action segmentation and detection, in: IEEE  
535 Conference on Computer Vision and Pattern Recognition, 2017, pp.  
536 156–165.
- 537 [11] B. Singh, T. K. Marks, M. Jones, O. Tuzel, M. Shao, A multi-stream bi-  
538 directional recurrent neural network for fine-grained action detection, in:  
539 IEEE Conference on Computer Vision and Pattern Recognition, 2016,  
540 pp. 1961–1970.

- 541 [12] D.-A. Huang, L. Fei-Fei, J. C. Niebles, Connectionist temporal model-  
542 ing for weakly supervised action labeling, in: European Conference on  
543 Computer Vision, 2016, pp. 137–153.
- 544 [13] Y. A. Farha, J. Gall, Ms-tcn: Multi-stage temporal convolutional net-  
545 work for action segmentation, arXiv preprint arXiv:1903.01945.
- 546 [14] S. Stein, S. J. McKenna, Combining embedded accelerometers with com-  
547 puter vision for recognizing food preparation activities, in: ACM In-  
548 ternational Joint Conference on Pervasive and Ubiquitous Computing,  
549 2013, pp. 729–738.
- 550 [15] A. Fathi, X. Ren, J. M. Rehg, Learning to recognize objects in egocen-  
551 tric activities, in: IEEE Conference on Computer Vision and Pattern  
552 Recognition, 2011, pp. 3281–3288.
- 553 [16] H. Kuehne, A. Arslan, T. Serre, The language of actions: Recovering  
554 the syntax and semantics of goal-directed human activities, in: IEEE  
555 Conference on Computer Vision and Pattern Recognition, 2014, pp.  
556 780–787.
- 557 [17] A. Fathi, J. M. Rehg, Modeling actions through state changes, in: IEEE  
558 Conference on Computer Vision and Pattern Recognition, 2013, pp.  
559 2579–2586.
- 560 [18] A. Fathi, A. Farhadi, J. M. Rehg, Understanding egocentric activities,  
561 in: IEEE International Conference on Computer Vision, 2011, pp. 407–  
562 414.
- 563 [19] S. Bhattacharya, M. M. Kalayeh, R. Sukthankar, M. Shah, Recognition  
564 of complex events: Exploiting temporal dynamics between underlying  
565 concepts, in: IEEE Conference on Computer Vision and Pattern Recog-  
566 nition, 2014, pp. 2235–2242.
- 567 [20] H. Kuehne, J. Gall, T. Serre, An end-to-end generative framework for  
568 video segmentation and recognition, in: IEEE Winter Conference on  
569 Applications of Computer Vision (WACV), 2016, pp. 1–8.
- 570 [21] K. Tang, L. Fei-Fei, D. Koller, Learning latent temporal structure for  
571 complex event detection, in: IEEE Conference on Computer Vision and  
572 Pattern Recognition, 2012, pp. 1250–1257.



- 573 [22] N. N. Vo, A. F. Bobick, From stochastic grammar to bayes network:  
574 Probabilistic parsing of complex activity, in: IEEE Conference on Com-  
575 puter Vision and Pattern Recognition, 2014, pp. 2641–2648.
- 576 [23] A. Richard, J. Gall, Temporal action detection using a statistical lan-  
577 guage model, in: IEEE Conference on Computer Vision and Pattern  
578 Recognition, 2016, pp. 3131–3140.
- 579 [24] A. Van Den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals,  
580 A. Graves, N. Kalchbrenner, A. W. Senior, K. Kavukcuoglu, Wavenet:  
581 A generative model for raw audio., ISCA Speech Synthesis Workshop  
582 (SSW) 125.
- 583 [25] L. Ding, C. Xu, Tricorner: A hybrid temporal convolutional and re-  
584 current network for video action segmentation, arXiv preprint arXiv:  
585 1705.07818.
- 586 [26] P. Lei, S. Todorovic, Temporal deformable residual networks for action  
587 segmentation in videos, in: IEEE Conference on Computer Vision and  
588 Pattern Recognition, 2018, pp. 6742–6751.
- 589 [27] P. Bojanowski, R. Lajugie, F. Bach, I. Laptev, J. Ponce, C. Schmid,  
590 J. Sivic, Weakly supervised action labeling in videos under ordering  
591 constraints, in: European Conference on Computer Vision, 2014, pp.  
592 628–643.
- 593 [28] L. Ding, C. Xu, Weakly-supervised action segmentation with iterative  
594 soft boundary assignment, in: IEEE Conference on Computer Vision  
595 and Pattern Recognition, 2018, pp. 6508–6516.
- 596 [29] H. Kuehne, A. Richard, J. Gall, Weakly supervised learning of actions  
597 from transcripts, *Computer Vision and Image Understanding* 163 (2017)  
598 78–89.
- 599 [30] A. Richard, H. Kuehne, J. Gall, Weakly supervised action learning with  
600 rnn based fine-to-coarse modeling, in: IEEE Conference on Computer  
601 Vision and Pattern Recognition, 2017, pp. 754–763.
- 602 [31] Z. Cai, N. Vasconcelos, Cascade r-cnn: Delving into high quality ob-  
603 ject detection, in: IEEE Conference on Computer Vision and Pattern  
604 Recognition, 2018, pp. 6154–6162.

- 605 [32] M. Amirul Islam, M. Rochan, N. D. Bruce, Y. Wang, Gated feedback  
606 refinement network for dense image labeling, in: IEEE Conference on  
607 Computer Vision and Pattern Recognition, 2017, pp. 3751–3759.
- 608 [33] G. Lin, A. Milan, C. Shen, I. Reid, Refinenet: Multi-path refinement  
609 networks for high-resolution semantic segmentation, in: IEEE confer-  
610 ence on computer vision and pattern recognition, 2017, pp. 1925–1934.
- 611 [34] H. Ding, X. Jiang, B. Shuai, A. Qun Liu, G. Wang, Context contrasted  
612 feature and gated multi-scale aggregation for scene segmentation, in:  
613 IEEE Conference on Computer Vision and Pattern Recognition, 2018,  
614 pp. 2393–2402.
- 615 [35] J. Gu, G. Wang, J. Cai, T. Chen, An empirical study of language cnn  
616 for image captioning, in: IEEE International Conference on Computer  
617 Vision, 2017, pp. 1222–1231.
- 618 [36] G. Varol, I. Laptev, C. Schmid, Long-term temporal convolutions for  
619 action recognition, IEEE Transactions on Pattern Analysis and Machine  
620 Intelligence 40 (6) (2017) 1510–1517.
- 621 [37] Y. N. Dauphin, A. Fan, M. Auli, D. Grangier, Language modeling with  
622 gated convolutional networks, in: International Conference on Machine  
623 Learning, 2017, pp. 933–941.
- 624 [38] M. Ranzato, S. Chopra, M. Auli, W. Zaremba, Sequence level training  
625 with recurrent neural networks, in: International Conference on Learn-  
626 ing Representations, 2016.
- 627 [39] R. J. Williams, Simple statistical gradient-following algorithms for con-  
628 nectionist reinforcement learning, Machine learning 8 (3-4) (1992) 229–  
629 256.
- 630 [40] C. Lea, A. Reiter, R. Vidal, G. D. Hager, Segmental spatiotemporal  
631 cnns for fine-grained action segmentation, in: European Conference on  
632 Computer Vision, 2016, pp. 36–52.
- 633 [41] C. Lea, R. Vidal, G. D. Hager, Learning convolutional action primitives  
634 for fine-grained action recognition, in: IEEE International Conference  
635 on Robotics and Automation, 2016, pp. 1642–1649.

- 636 [42] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recog-  
637 nition, in: IEEE Conference on Computer Vision and Pattern Recogni-  
638 tion, 2016, pp. 770–778.