

# Research Report

## 视频场景文字检测与识别算法研究

2021年 10月

指导老师：刘绍辉 汇报人：研一 舒言

# 目录

## CONTENTS

- ▷ 第一部分 『任务概述』
- ▷ 第二部分 『主流方法』
- ▷ 第三部分 『未来展望』



# 第一部分

## 任务概述



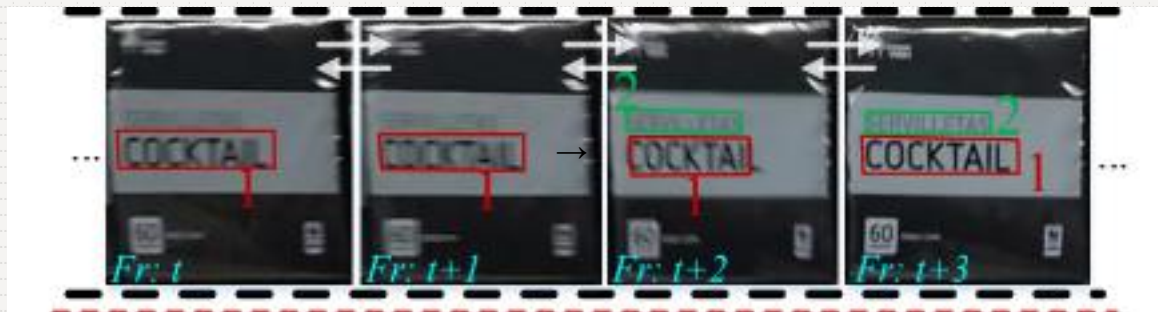
## ICDAR 2021 Competition on Scene Video Text Spotting

Zhanzhan Cheng<sup>1,2\*</sup>, Jing Lu<sup>2\*</sup>, Baorui Zou<sup>3\*</sup>, Shuigeng Zhou<sup>3</sup>, and Fei Wu<sup>1</sup>

<sup>1</sup> Zhejiang University, Hangzhou, China  
{11821104,wufei}@zju.edu.cn

<sup>2</sup> Hikvision Research Institute, Hangzhou, China  
lujing6@hikvision.com

<sup>3</sup> Fudan University, Shanghai, China  
{18210240270,sgzhou}@fudan.edu.cn



Text Detection

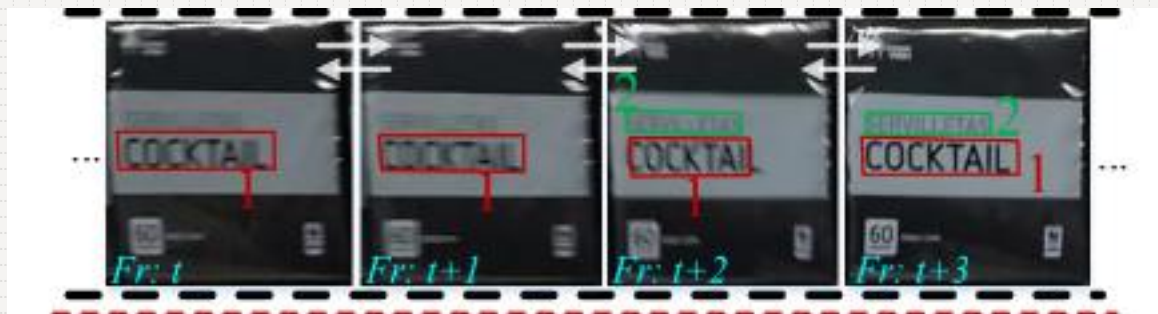


Text Tracking



Text Recognition





Text Detection

Text Tracking

Text Recognition

```

{
  "video_19_4/67.jpg": {
    "content_ann": {
      "bboxes": [
        [210, 646, 746, 670],
        [763, 752, 213, 213],
        [704, 513, 355, 677],
        [275, 692, 327, 517],
        [397, 765, 619, 773],
        [634, 764, 658, 739],
        [622, 622]
      ]
    }
  },
  "video_19_6/405.jpg": {
    "content_ann": {
      "bboxes": [
        [765, 619, 773, 634],
        [764, 658, 739, 622]
      ]
    }
  }
}

```

(a)

```

{
  "video_10_3": {
    "1304": {
      "tracks": [
        "158,938_569_962_569_962_589_938_589"
      ]
    },
    "1348": {
      "tracks": [
        "310,605_779_648_779_648_811_605_811",
        "313,586_799_632_799_632_827_586_827",
        "309,607_782_647_782_647_810_607_810",
        "311,604_783_639_783_639_808_604_808",
        "312,599_791_639_791_639_819_599_819"
      ]
    }
  },
  "video_9_1": {
    "1140": {
      "tracks": [
        "40,181_1102_520_1102_520_1181_181_1181",
        "30,526_1156_871_1156_871_1234_526_1234",
        "41,145_1099_484_1099_484_1183_145_1183",
        "36,397_1176_735_1176_735_1257_397_1257",
        "45,124_1172_446_1172_446_1259_124_1259",
        "44,104_1151_438_1151_438_1235_104_1235"
      ]
    }
  }
}

```

(b)

```

{
  "video_10_3": {
    "1304": {
      "tracks": [
        "158,938_569_962_569_962_589_938_589,DAVAR"
      ],
      "text": "DAVAR"
    },
    "1348": {
      "tracks": [
        "310,605_779_648_779_648_811_605_811,LAB",
        "313,586_799_632_799_632_827_586_827,LAB",
        "309,607_782_647_782_647_810_607_810,LAB",
        "311,604_783_639_783_639_808_604_808,LAB",
        "312,599_791_639_791_639_819_599_819,LAB"
      ],
      "text": "LAB"
    }
  },
  "video_9_1": {
    "1140": {
      "tracks": [
        "40,181_1102_520_1102_520_1181_181_1181,JUICE",
        "30,526_1156_871_1156_871_1234_526_1234,JUICE",
        "41,145_1099_484_1099_484_1183_145_1183,JUICE",
        "36,397_1176_735_1176_735_1257_397_1257,JUICE",
        "45,124_1172_446_1172_446_1259_124_1259,JUICE",
        "44,104_1151_438_1151_438_1235_104_1235,JUICE"
      ],
      "text": "JUICE"
    }
  }
}

```

(c)



(1) outdoor shopping mall



(2) pedestrian



(3) fingerpost



(4) book opening



(5) digital screen



(6) indoor shopping mall



(7) inside shops



(8) supermarket



(9) metro station



(10) restaurant



(11) office building



(12) hotel



(13) bus/railway station



(14) bookstore



(15) street view



(16) inside train



(17) train watch



(21) shopping bags



(18) city road



(19) harbor surveillance

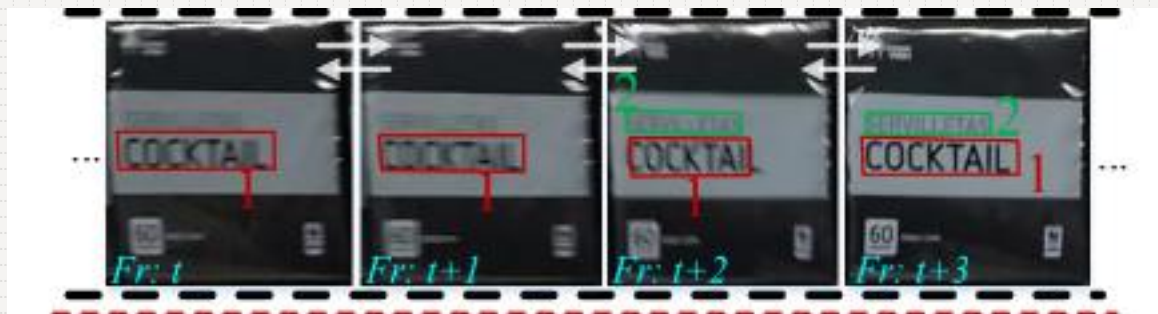


(20) highway

blurring, rotation,  
illumination.....

## 1

## 评测指标



Text Detection

Text Tracking

Text Spotting

IoU-based

Precision

Recall

F-measure

Similar to multi-object  
tracking

ATA

MOTA

MOTP

Detection+Tracking  
+recognition

Precision

Recall

F-measure

$$MOTA = 1 - \frac{\sum(FN + FP + IDSW)}{\sum GT} \in (-\infty, 1]$$



User ID	Rank	F-score <sub>s</sub>	Precision <sub>s</sub>	Recall <sub>s</sub>	ATA <sub>s</sub>	MOTA <sub>s</sub>	MOTP <sub>s</sub>	Affiliations
tianqihenhao	1	0.5308	0.6655	0.4414	0.4549	0.5913	0.8421	TEG, Tencent
DXM-DI-AI -CV-TEAM	2	0.4755	0.6435	0.3770	0.4188	0.4960	0.8142	DuXiaoman Financial
panda12	3	0.4183	0.5243	0.3479	0.3579	0.5179	0.8427	IA, CAS
lzneu09	4	0.3007	0.3611	0.2576	0.2737	0.4255	0.8330	Northeastern University
yucheng3	5	0.2964	0.3506	0.2567	0.2711	0.4246	0.8332	University of Chinese Academy of Sciences
tangyejun	6	0.2284	0.2527	0.2084	0.2121	0.3676	0.8337	*
tiendv	7	0.0813	0.1402	0.0572	0.0802	0.0887	0.7976	University of Information Technology
enderloong	8	0.0307	0.0239	0.0429	0.0357	0.0159	0.7813	*
colorr	9	0.0158	0.0085	0.1225	0.0146	0.0765	0.8498	*
weijiawu3	10	0.0077	0.0041	0.0550	0.0088	-0.1530	0.7670	Zhejiang University
BOE-AIoT-CTO	11	0.0000	0.0000	0.0000	0.0000	-0.0003	0.0000	BOE



第二部分



主流方法



## You Only Recognize Once: Towards Fast Video Text Spotting

Zhanzhan Cheng<sup>12\*</sup>, Jing Lu<sup>2\*</sup>, Yi Niu<sup>2</sup>, Shiliang Pu<sup>2</sup>, Fei Wu<sup>1+</sup>, Shuigeng Zhou<sup>3</sup>

<sup>1</sup>Zhejiang University, Hangzhou, China

chengzhanzhan@hikvision.com, wufei@cs.zju.edu.cn

<sup>2</sup>Hikvision Research Institute, China

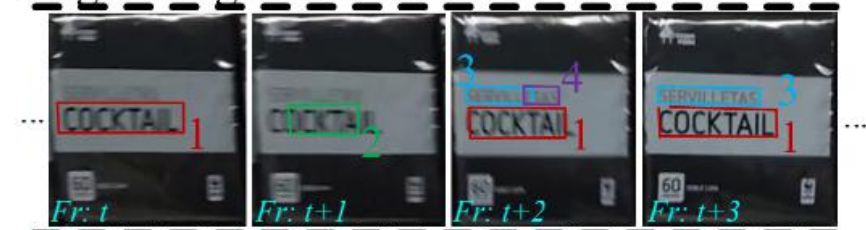
lujing6, niuyi, pushiliang@hikvision.com

<sup>3</sup>Fudan University, Shanghai, China

sgzhou@fudan.edu.cn



## Single Image Text Detection



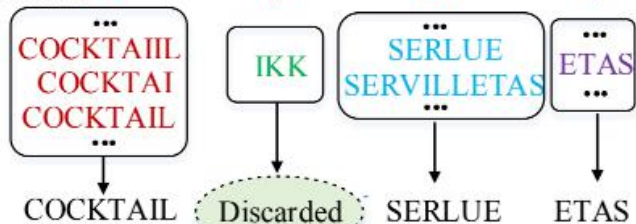
## Frame-wise Recognition

$$t \begin{cases} \text{COCKTAIL} \\ \text{COCKTAIL} \\ \text{COCKTAIL} \end{cases} \quad t+1 \begin{cases} \text{IKK} \end{cases} \quad t+2 \begin{cases} \text{COCKTAIL} \\ \text{SERLUE} \\ \text{ETAS} \end{cases} \quad t+3 \begin{cases} \text{COCKTAIL} \\ \text{SERVILLETAS} \end{cases}$$

## Tracking

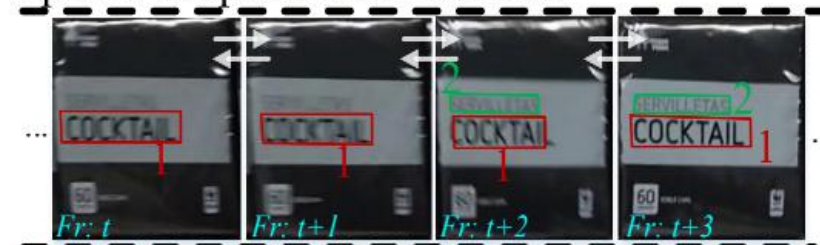


## Post-process

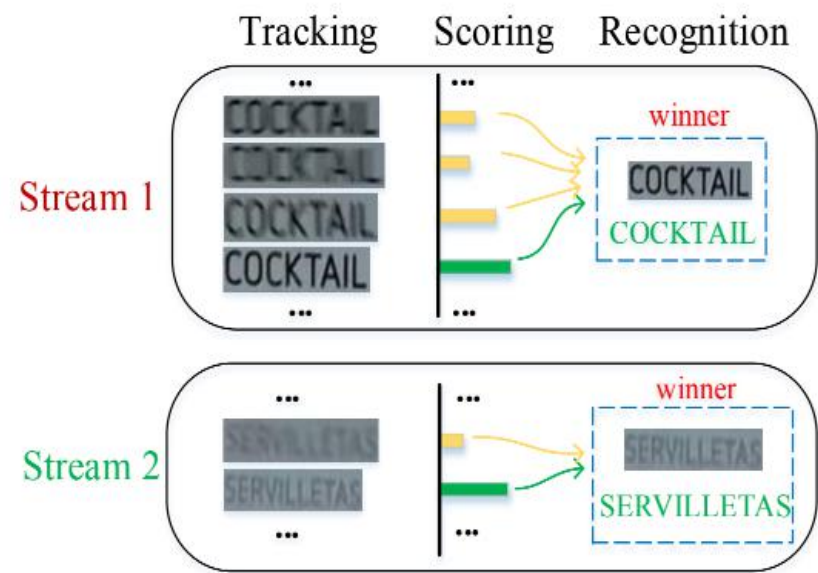


(a)

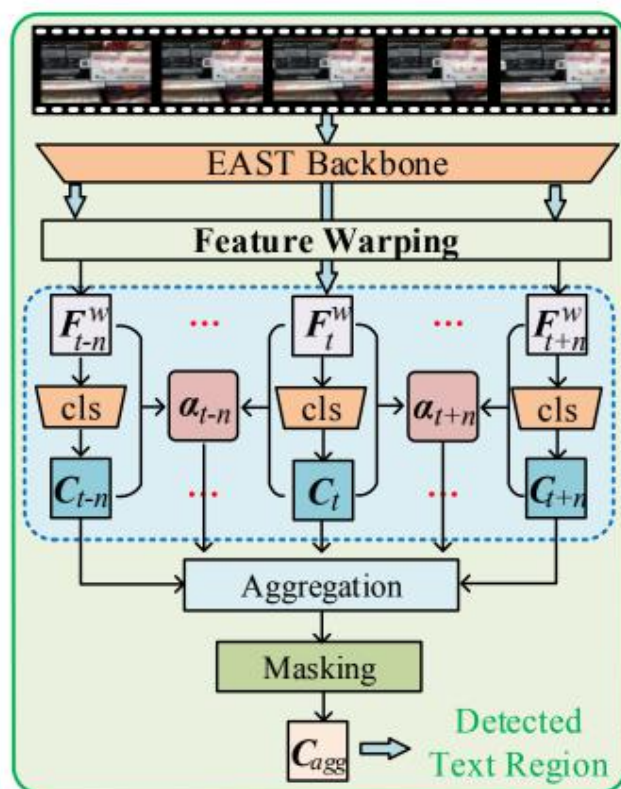
## Spatial-temporal Text Detector



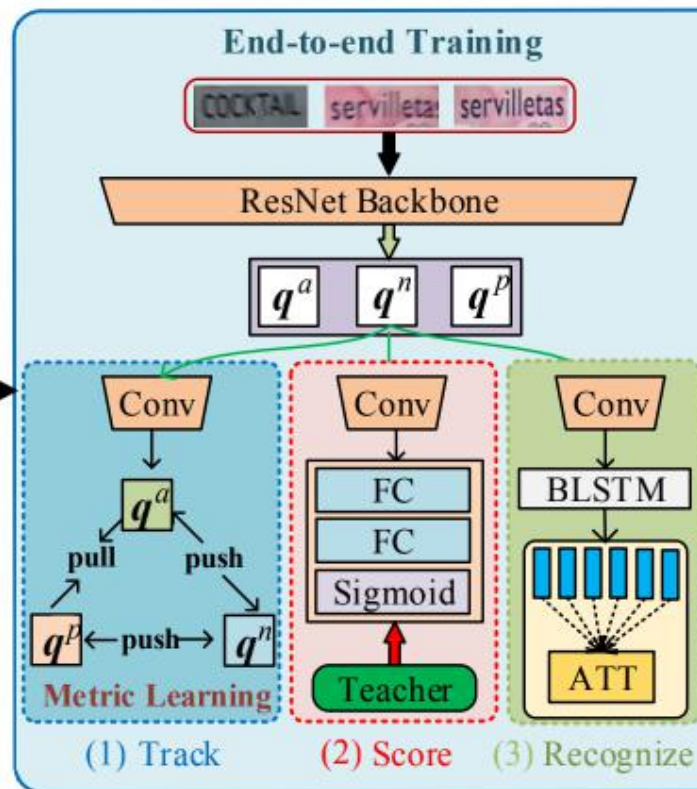
## Text Recommender



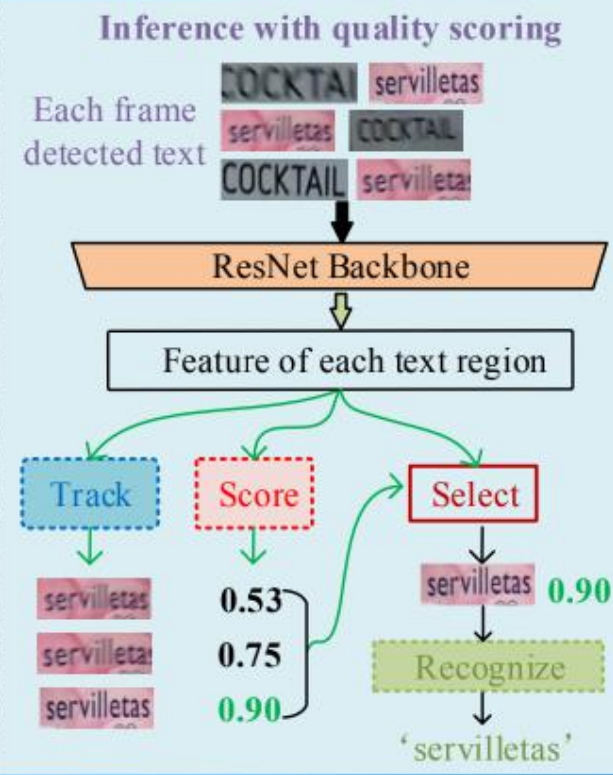
(b)



(a) Spatial-temporal Detector



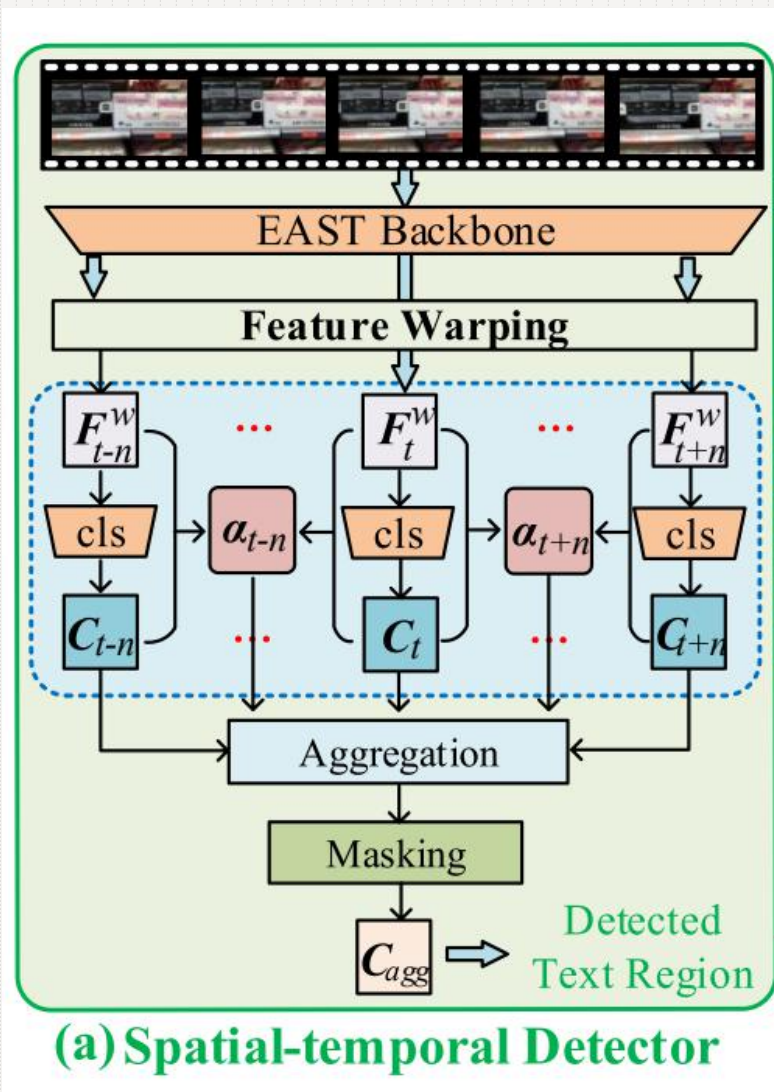
(b) Text Recommender





## 2

## Detection Module



Input

 $I_{t-n} \dots I_t \dots I_{t+n}$ 

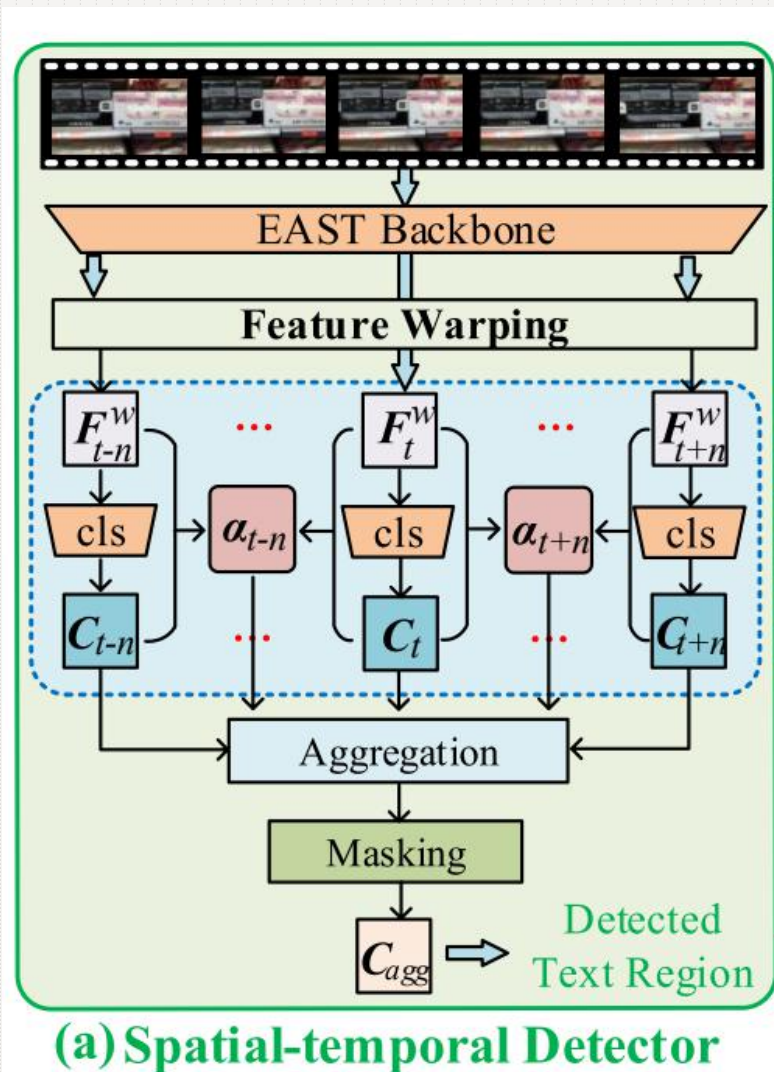
Model

Output

Text region mask

## 2

## Detection Module


 $I_{t-n} \dots I_t \dots I_{t+n}$ 
 $F_{t-n} \dots F_t \dots F_{t+n}$ 

$$F_{t+i}^w = \text{Warp}(F_{t+i}, \text{flow}_{(t+i,t)}),$$

MLP-based Classifier

 $C_{t-n} \dots C_t \dots C_{t+n}$ 

$$F_{t+i}^{\text{trans}} = \text{ReLU}(\text{BN}(W F_{t+i}^w + b)),$$

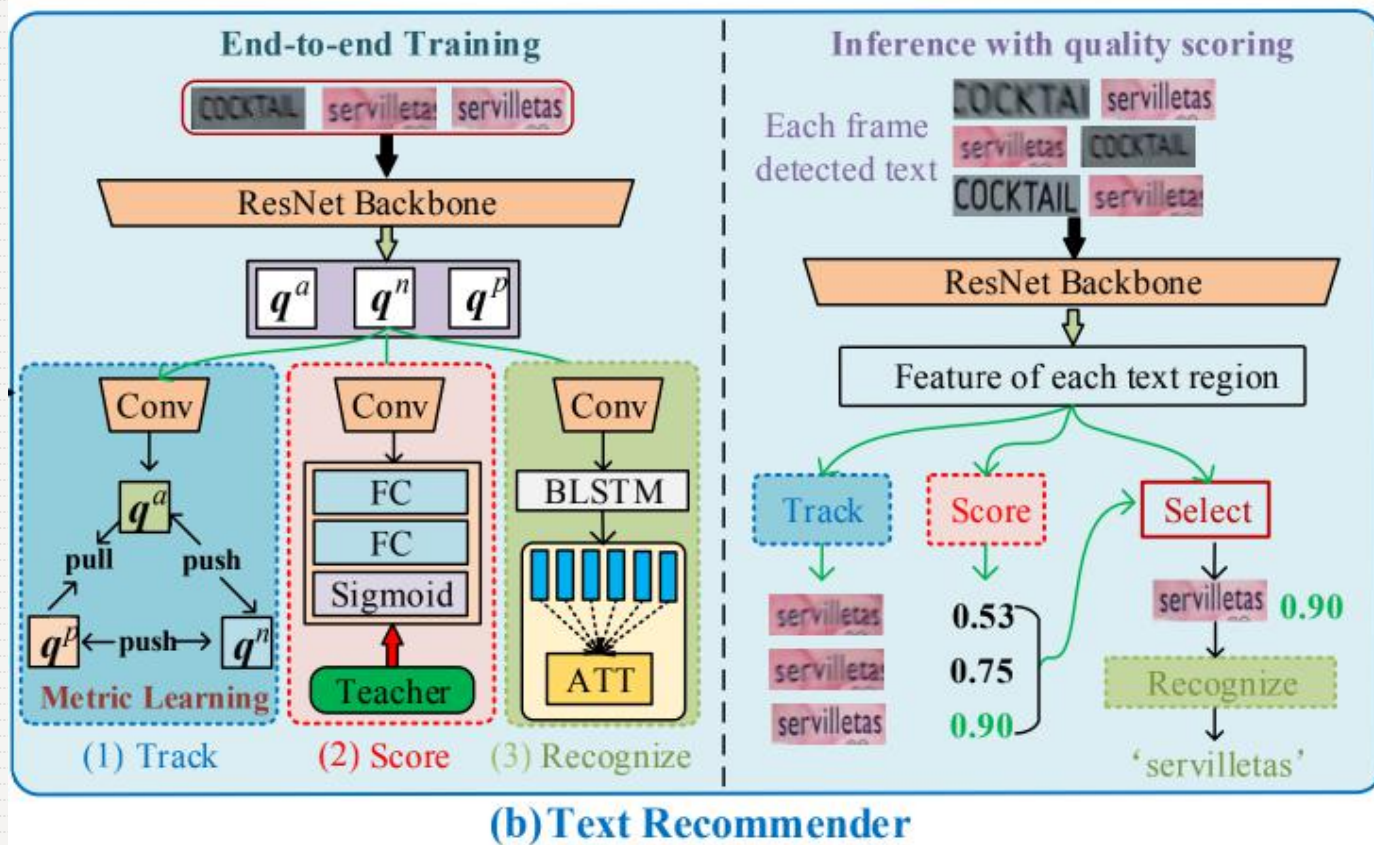
$$\text{Sim}_{t+i,t} = F_{t+i}^{\text{trans}} \odot F_t^{\text{trans}}$$

$$a_{t+i} = \frac{\exp(\text{Sim}_{t+i,t} \odot C_{t+i})}{\sum_{i'=-n}^n \exp(\text{Sim}_{t+i',t} \odot C_{t+i'})}.$$

$$C_{t,\text{agg}} = \sum_{i=-n}^n a_{t+i} * C_{t+i}$$

Detected  
Text Region





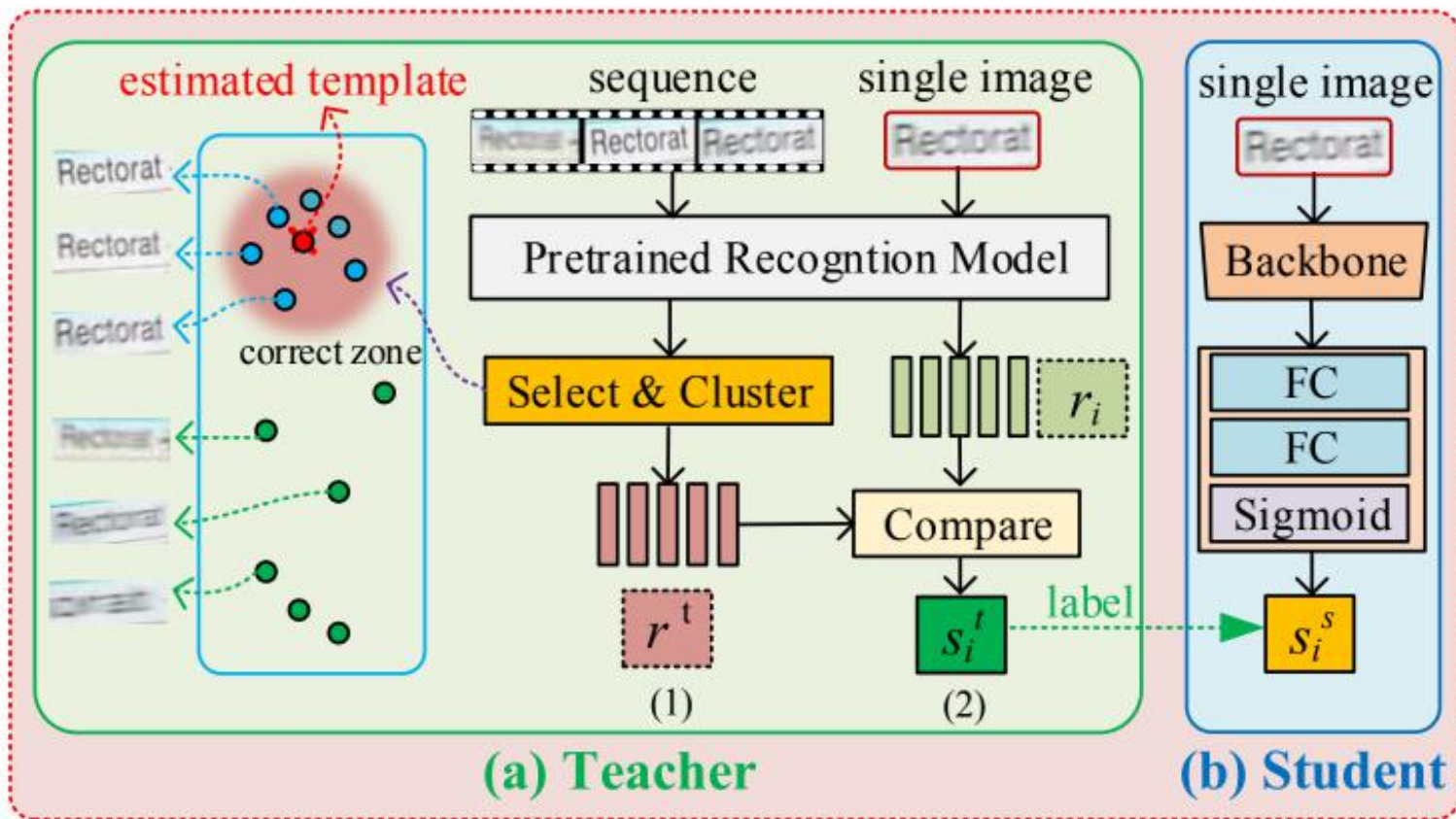
Quality Score

Tracking

Recognizing

## 2

## Text Recommender--Quality score

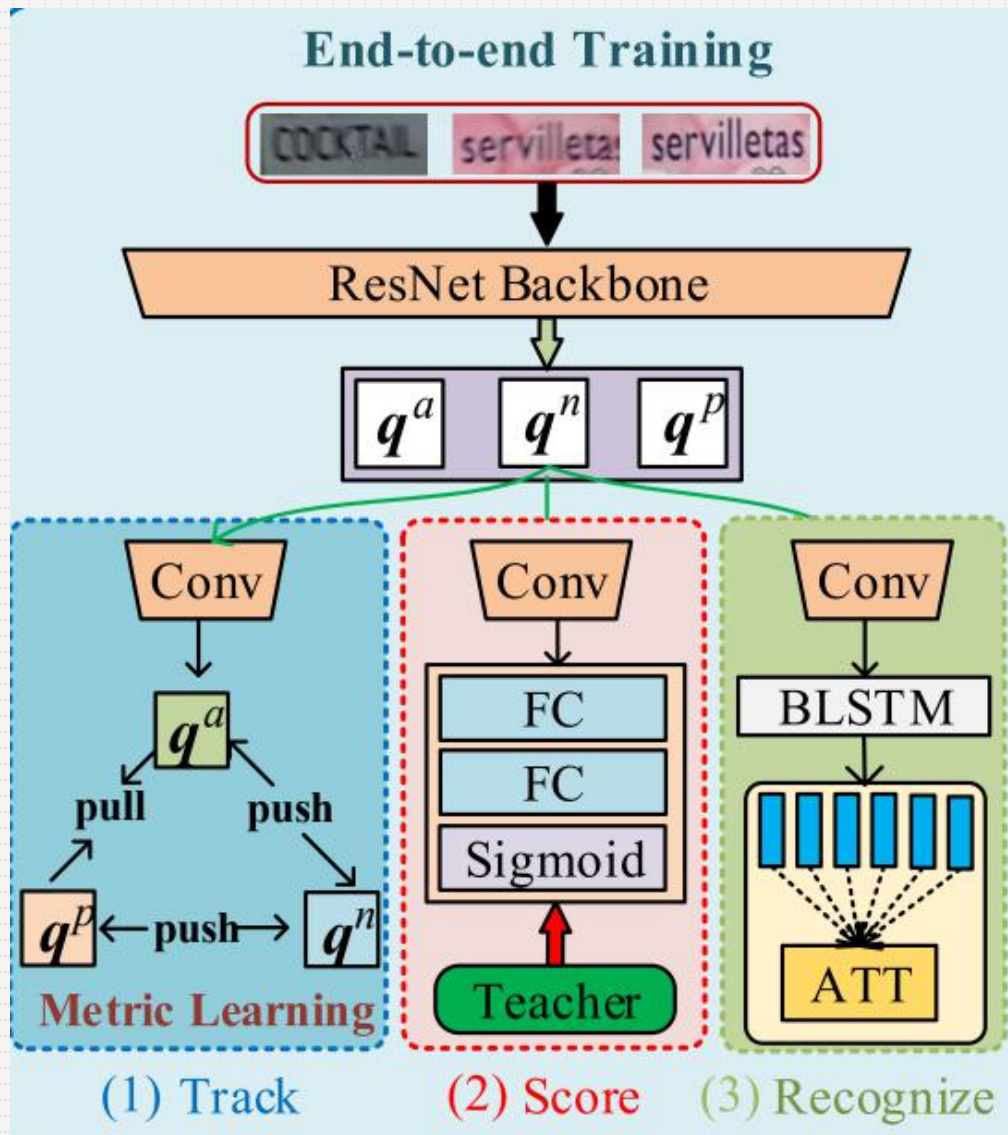


$$r^t = kmeans(r_1^{cor}, r_1^{cor}, \dots, r_k^{cor}),$$

$$s_i^t = \frac{r^t \odot r_i}{||r^t|| * ||r_i||},$$

$$\mathcal{L}_S = \frac{1}{N} \sum_{i=0}^N ||s_i^t - s_i^s||$$





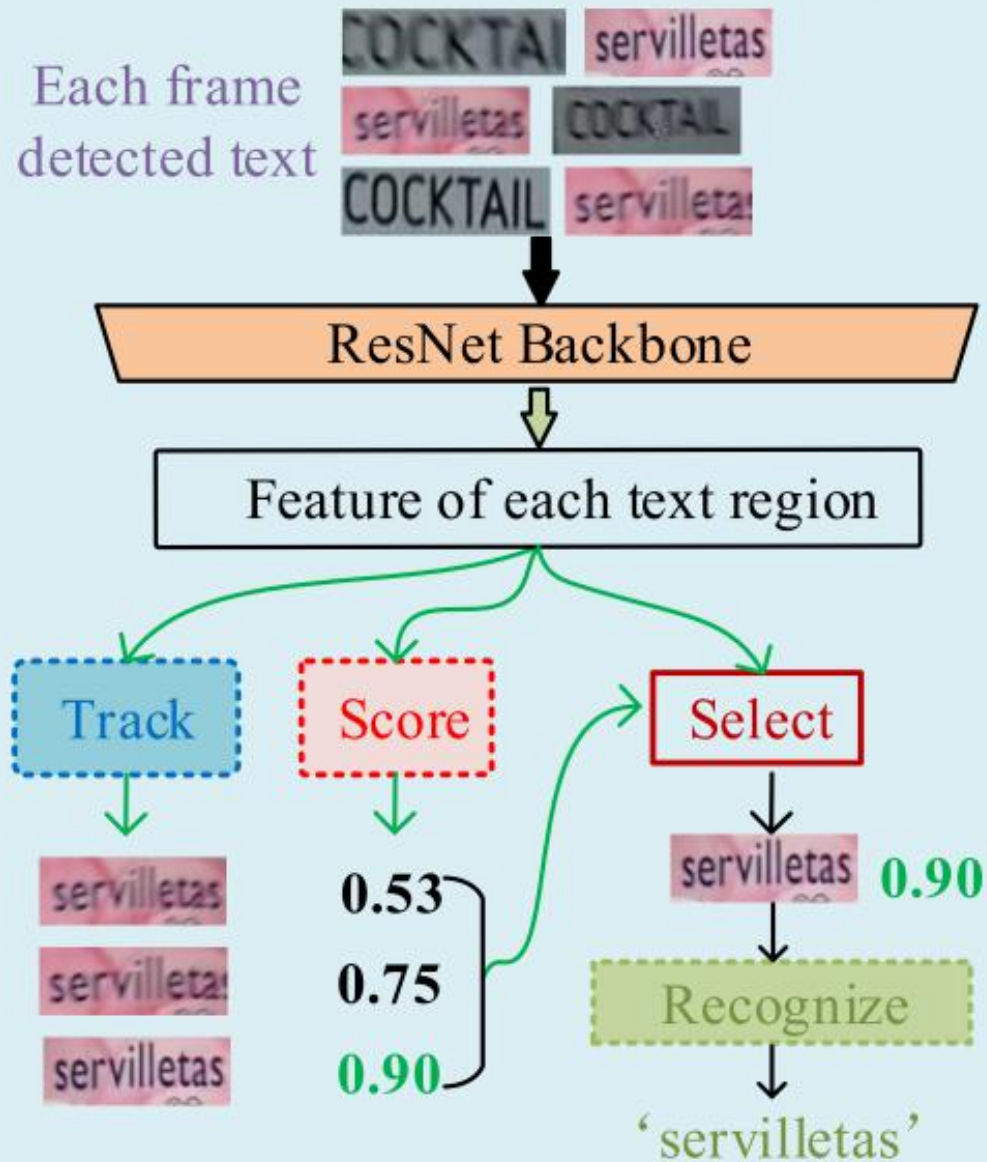
$$\mathcal{L}_T = \mathcal{L}_{contra} + \lambda_t \mathcal{L}_{triplet},$$

$$\mathcal{L} = \lambda_1 \mathcal{L}_T + \lambda_2 \mathcal{L}_S + \lambda_3 \mathcal{L}_R,$$

## 2

## Inference

## Inference with quality scoring





## 2

## Ablation Experiments

Methods	$QSHR$ (IC13/IC15)	$RCR$ (IC13/IC15)	$FPS$
PCW	74.55/75.83	66.06/66.32	4.52
HFP	75.32/76.34	68.30/68.56	
TR ( $\mathcal{L}_S$ )	77.89/79.69	68.89/69.41	324.58
TR ( $\mathcal{L}_S + \mathcal{L}_T$ )	78.64/80.36	69.12/69.82	
TR ( $\mathcal{L}_S + \mathcal{L}_R$ )	81.23/83.03	69.92/70.69	
TR ( $\mathcal{L}$ )	<b>81.74/83.29</b>	<b>70.18/70.95</b>	

Methods	$QSHR$ (IC13/IC15)	$RCR$ (IC13/IC15)
PCW	41.73/45.66	59.78/60.62
HFP	39.37/41.73	58.96/60.06
TR	<b>51.18/54.33</b>	<b>66.14/67.71</b>

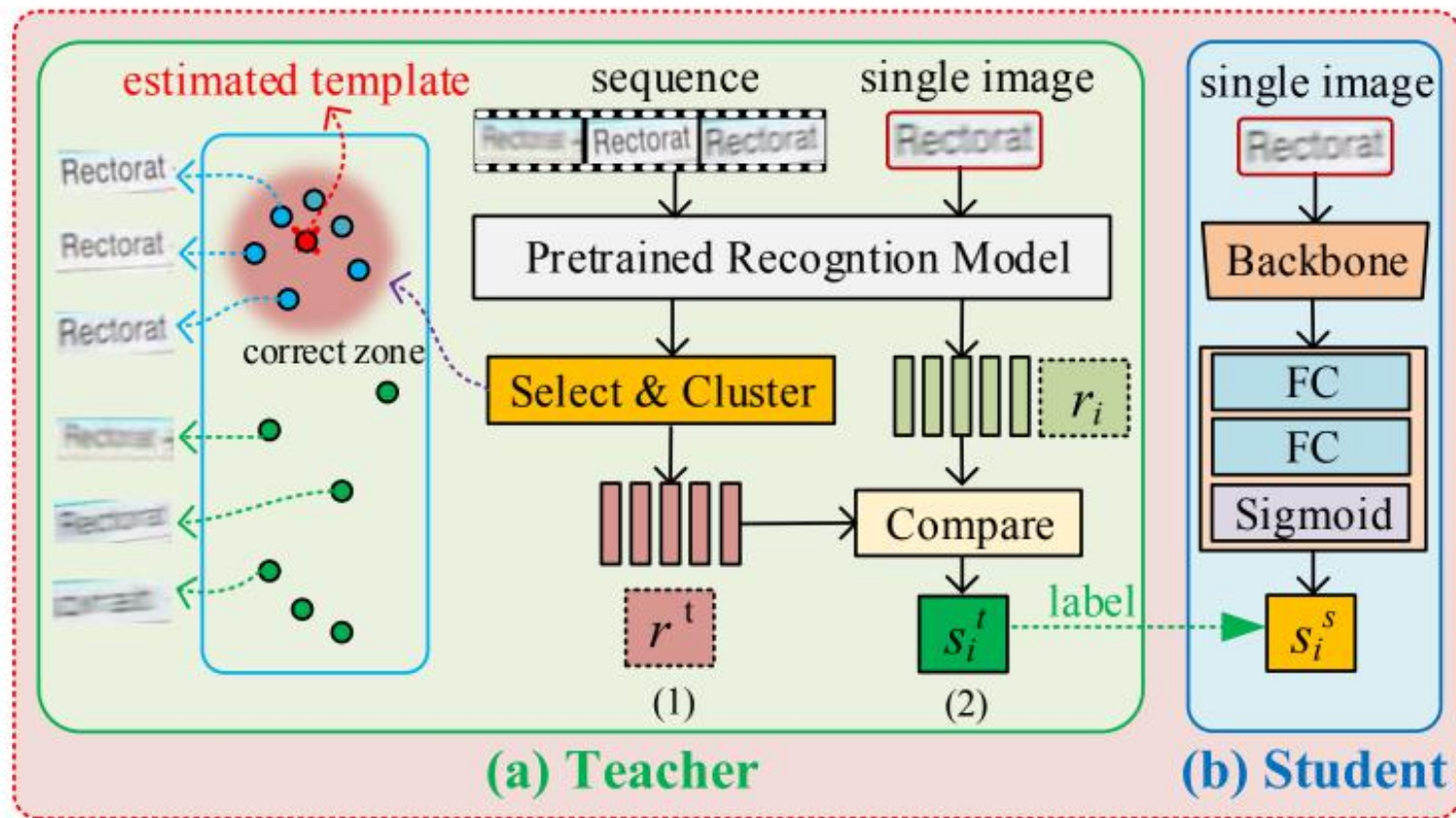
## 第三部分

# 未来展望













# 感谢聆听！

THANK YOU FOR WATCHING!

---

2021年9月

---

