# **MEAN: Multi-Element Attention Network for Scene Text Recognition**

Ruijie Yan\*, Liangrui Peng\*, Shanyu Xiao\*, Gang Yao\*, and Jaesik Min<sup>†</sup> \*Beijing National Research Center for Information Science and Technology Dept. of Electronic Engineering, Tsinghua University, Beijing, China *†* Hyundai Motor Group AIRS Company, Seoul, Korea

#### Introduction

We propose a novel multi-element attention (MEA) mechanism to exploit geometric structures from local to global levels in feature maps extracted from a scene text image. A multi-element attention network (MEAN) is constructed by using three types of MEAs. Experiments on English and Chinese scene text recognition experiments demonstrates the effectiveness of MEAN on regular, irregular, and multi-oriented texts.



System framework of MEAN that consists of a CNN, an encoder, and a decoder.

## Method

### **Multi-Element Attention (MEA)**

- Idea
  - Incorporating graph structure modeling into self-attention mechanism and assigning various adjacency matrices to the graph.

### Experiment

**English scene text recognition** 

Word Accuracy (%) across different models and datasets

Model	IIIT5k	SVT	<b>IC03</b>	<b>IC13</b>	<b>IC15</b>	SVTP	CUTE
Mask TextSpotter	95.3	91.8	95.0	95.3	78.2	83.6	88.5
SAR	95.0	91.2	-	94.0	78.8	86.4	89.6
ASTER	93.4	89.5	94.5	91.8	76.1	78.5	79.5



- Three implementations of  $AXW_Q$  and  $BXW_K$ 
  - **MEA-Local**: 1 × 1 convolutions
  - **MEA-Neighbor**:  $m \times n$  convolutions
  - **MEA-Global**: graph convolutions

SRN	94.8	91.5	-	95.5	82.7	85.1	87.8
MEAN	95.9	94.3	95.9	95.1	79.7	86.8	87.2

• Multi-oriented Chinese scene text recognition

Word Accuracy (%) of different models on multi-oriented texts

Tost sot	J	Baselin	e		MEAN	J
	Η	V	H & V	Η	V	H & V
Horizontal	74.2	-	52.2	77.4	-	81.6
Vertical	_	74.6	36.0	_	78.6	86.0

#### **Visualization of attention scores**



#### **Multi-Element Attention Network (MEAN)**

#### • Framework

- MEAN consists of a U-shaped CNN, an encoder with three MEAs, and a decoder.
- **Orientational positional encoding**  $\bullet$

$$PE_{(i,j,2k)}^{H} = \sin(j/L^{2k/d}) \qquad PE_{(i,j,2k)}^{V} = \sin(i/L^{2k/d})$$
$$PE_{(i,j,2k+1)}^{H} = \cos(j/L^{2k/d}) \qquad PE_{(i,j,2k+1)}^{V} = \cos(i/L^{2k/d})$$

**Recognition Examples** 







鸟语花香 Output: 鸟语花香

starbucks GT: Output: starbucks

(b) Skewed texts

(c) Multi-oriented texts

#### **Acknowledgment:**

- This research is supported by a joint research project between Hyundai Motor Group AIRS Company and Tsinghua University.
- The second author is partially supported by National Key R&D Program of China and a grant from the Institute for Guo Qiang, Tsinghua University.