



# Industrial Scene Text Detection with Refined Feature-attentive Network

Tongkun Guan<sup>1</sup>, Chaochen Gu<sup>1</sup>, Changsheng Lu<sup>2</sup>, Jingzheng Tu<sup>1</sup>, Qi Feng<sup>1</sup>, Kaijie Wu<sup>1</sup>, Xinping Guan<sup>1</sup>

<sup>1</sup> Shanghai Jiao Tong University, <sup>2</sup> The Australian National University

## Problem Definition and Contribution

**Goal:** Detecting the multi-oriented texts on the surface of metal parts in industrial scenarios.

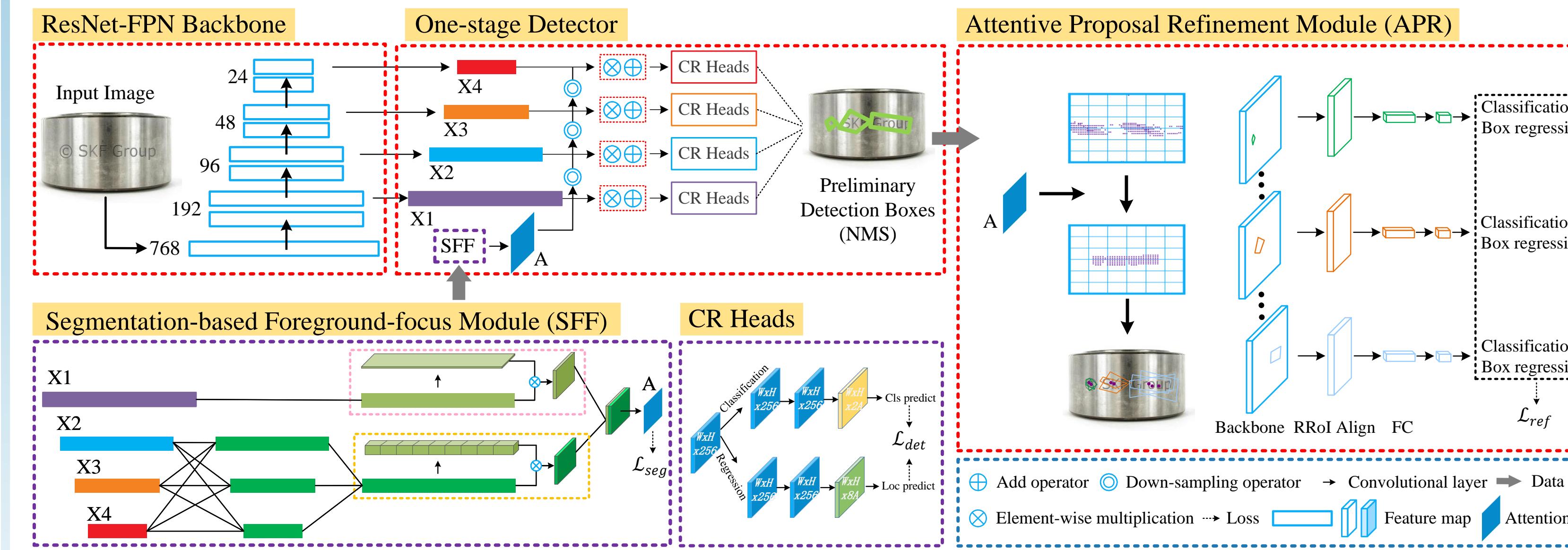


**Key Contributions:** A flexible deep learning framework for industrial text detection that

- focuses on more detailed text foreground features to weaken the influence of other factors (e.g., low visual contrast, corroded surfaces, and complex backgrounds).
- excavates high-quality foreground boxes for localization correction.
- achieves state-of-the-art performance on the industrial dataset and robustly detects horizontal texts, multi-oriented texts, and multi-language texts of natural scenes.
- establishes the first benchmark dataset (Metal Part Surface Character Dataset, MPSC) to promote in-depth research on industrial text detection.

## Method

**Network Architecture:** RFN consists of three components, namely a ResNet-FPN backbone, a one-stage detector, and an attentive proposal refinement module.



**Post-processing:** We combine the instance score  $S_I$  and classification score  $S_c$  to generate a new confidence score. And then, they are fed into the NMS algorithm to get the best prediction boxes:

$$S_I = \frac{\sum_{j=1}^N \rho_j}{N}, S' = e^{S_c} \left(1 + \mu \frac{e^{S_I}}{e^{1-S_I}}\right) \quad (1)$$

## Optimization

**Assumption:**  $\omega_i$  and  $\omega_i^*$  are the confidence score of pixel  $i$  in  $S_{gt}$  and the attention map  $A$ .

**Loss function:**

$$\omega_d = \omega_i - \omega_i^*, \mathcal{D}_a = \frac{\sum_{i=1}^N 1_{[\omega_d \geq \frac{1}{2}]} (1 - \omega_i^*)}{\sum_{i=1}^N (\omega_i^*)}, \mathcal{D}_b = \frac{\sum_{i=1}^N 1_{[-\omega_d \geq \frac{1}{2}]} \omega_i^*}{\sum_{i=1}^N (\omega_i^*)}, \quad (2)$$

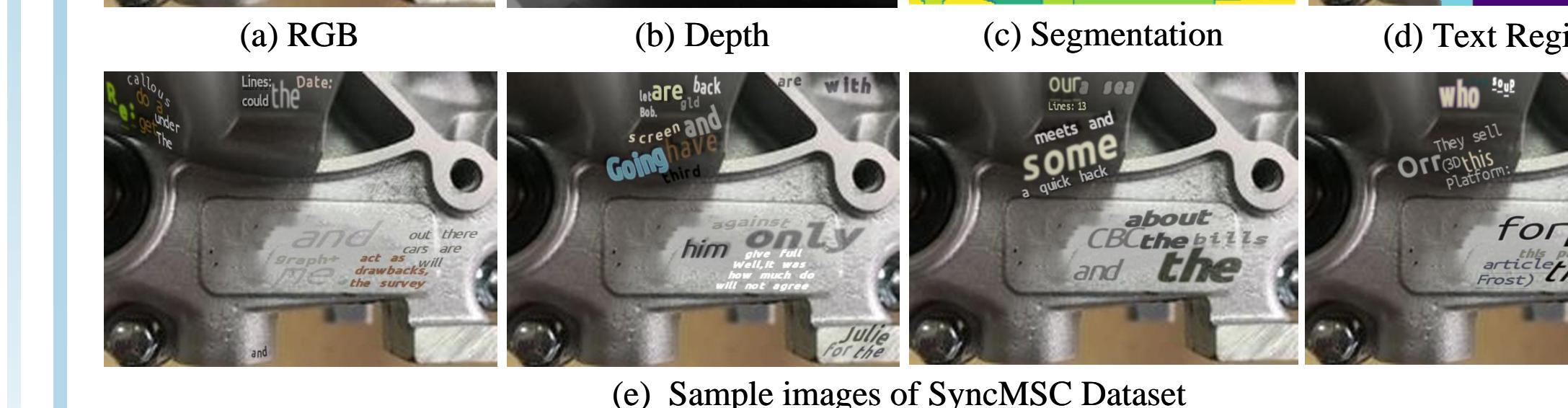
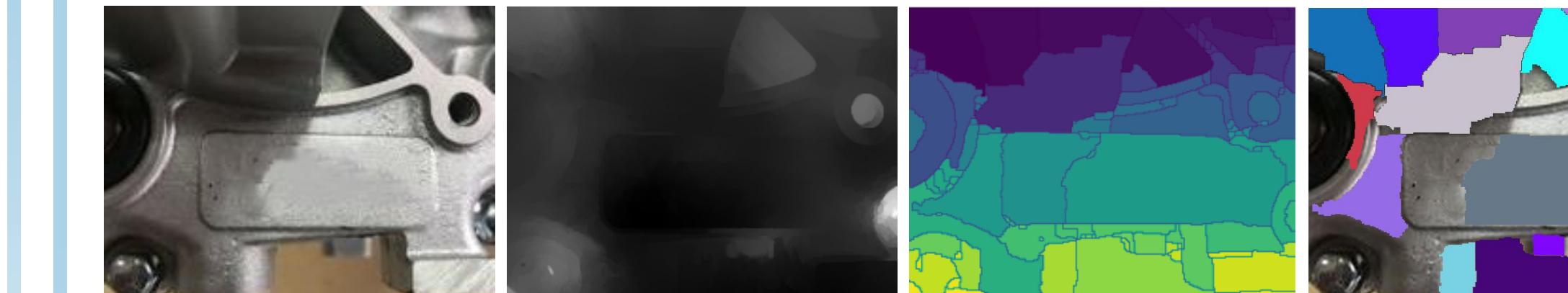
$$\mathcal{L}_g = \begin{cases} \mathcal{D}_a & \text{if } \mathcal{D}_b < \Delta, \\ \mathcal{D}_a + \mathcal{D}_b - \Delta & \text{if } \mathcal{D}_b \geq \Delta \end{cases}, \mathcal{L}_{seg} = \mathcal{L}_d + e^{-1 * \mathcal{L}_d * \gamma} * \mathcal{L}_g,$$

where  $\mathcal{L}_d$  denotes the dice loss,  $N$  is the number of pixels in the attention map  $A$ .

**Main Idea:** We allow certain false-positive classification results in exchange for detecting more textual features in low-contrast and indistinguishable regions when the  $\mathcal{L}_d$  stabilizes.

## Experiments & Results

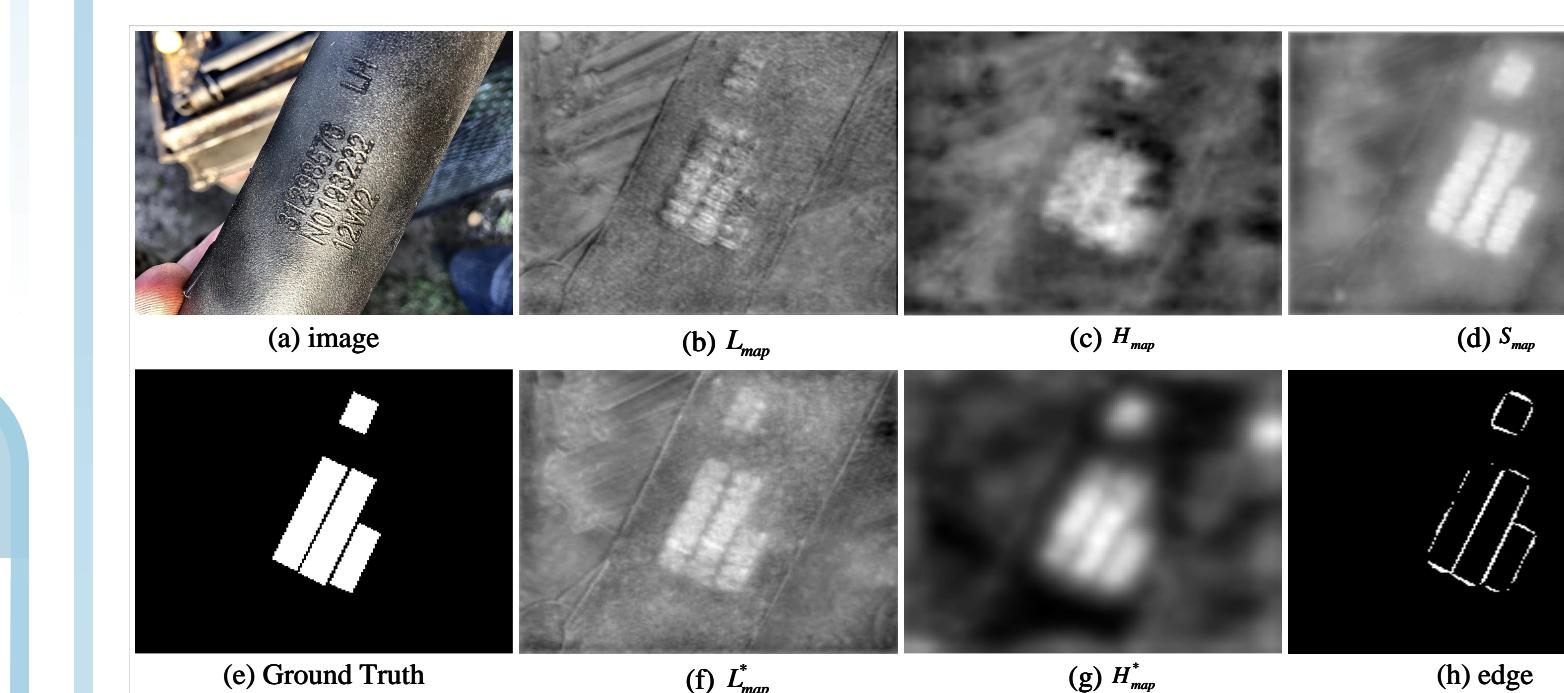
**Synthetic Datasets (SynthMPSC, 98962 images):**



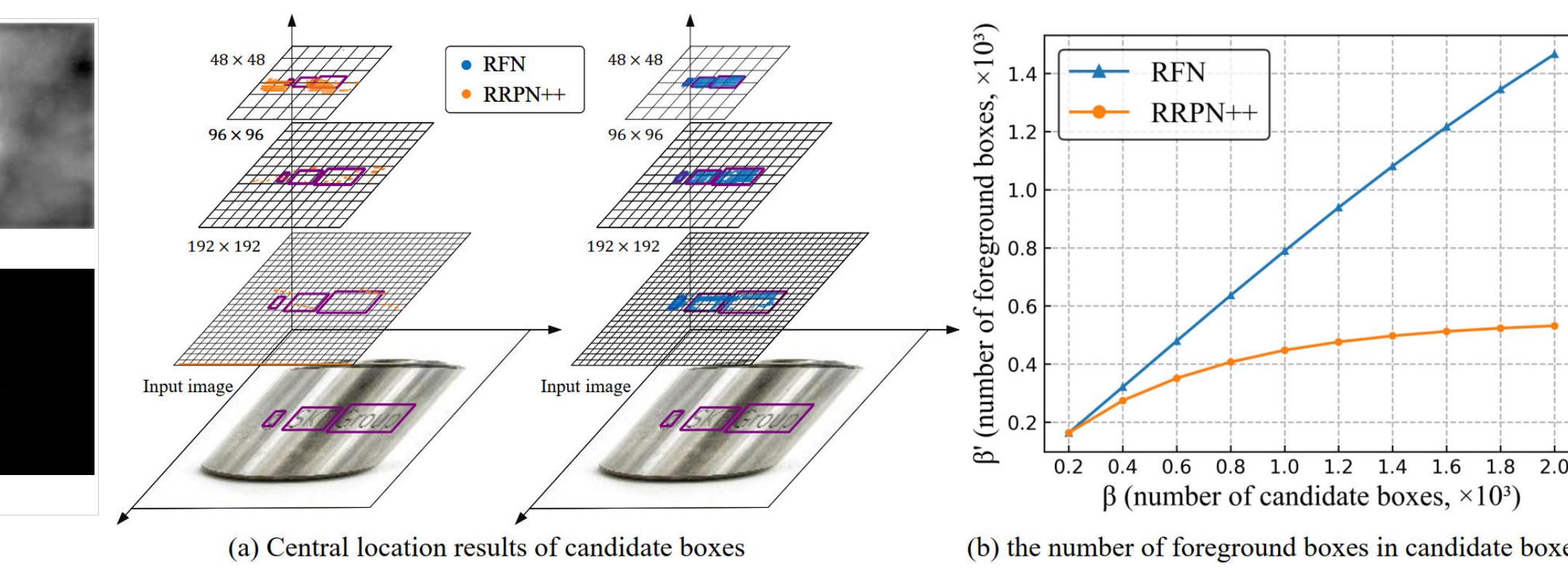
**Self-built Benchmark Dataset (MPSC, 3194 images):**



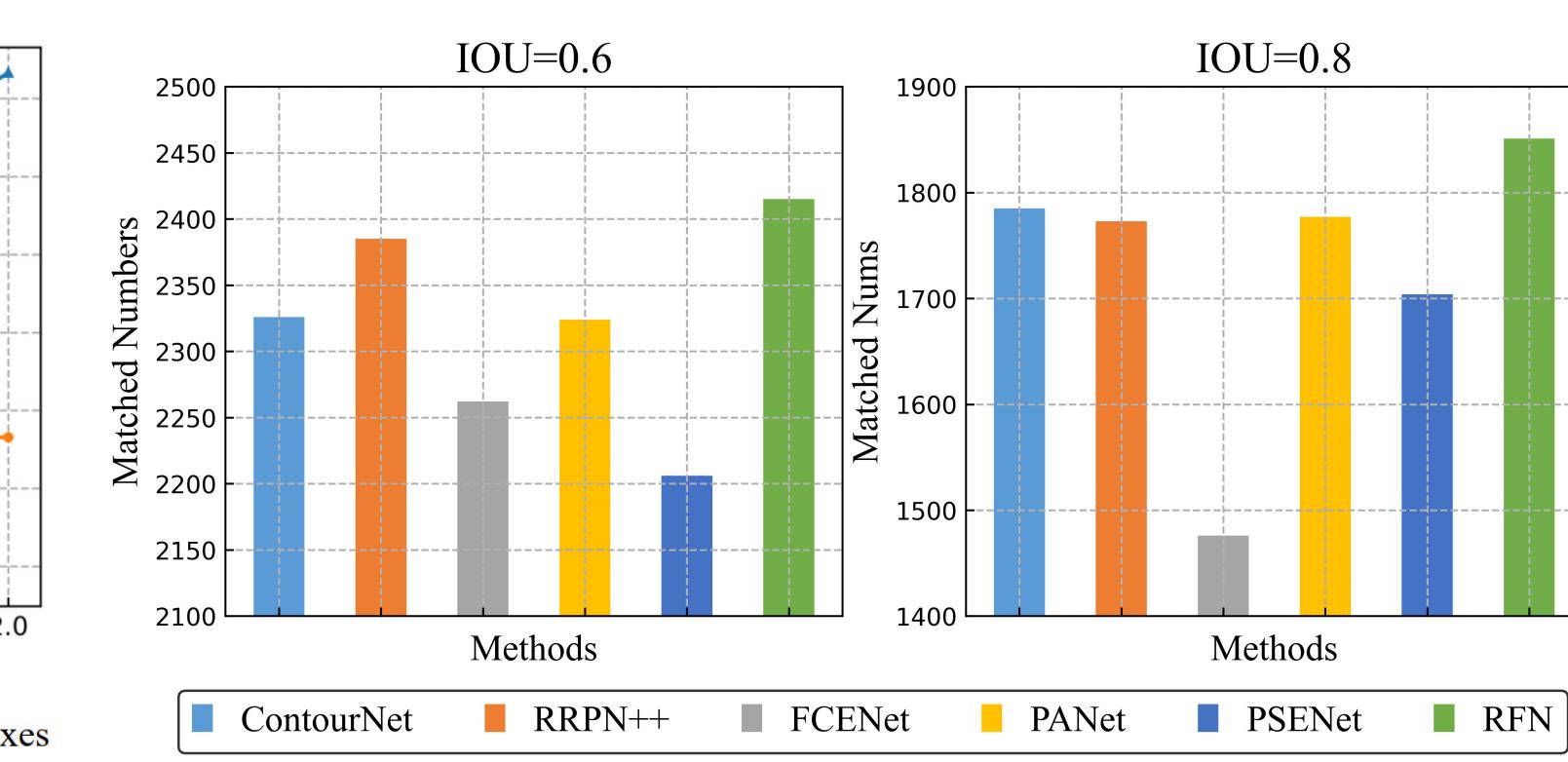
**Feature Visualization (SFF):**



**Central location results of candidate boxes:**



**Location Correction Results (APR):**

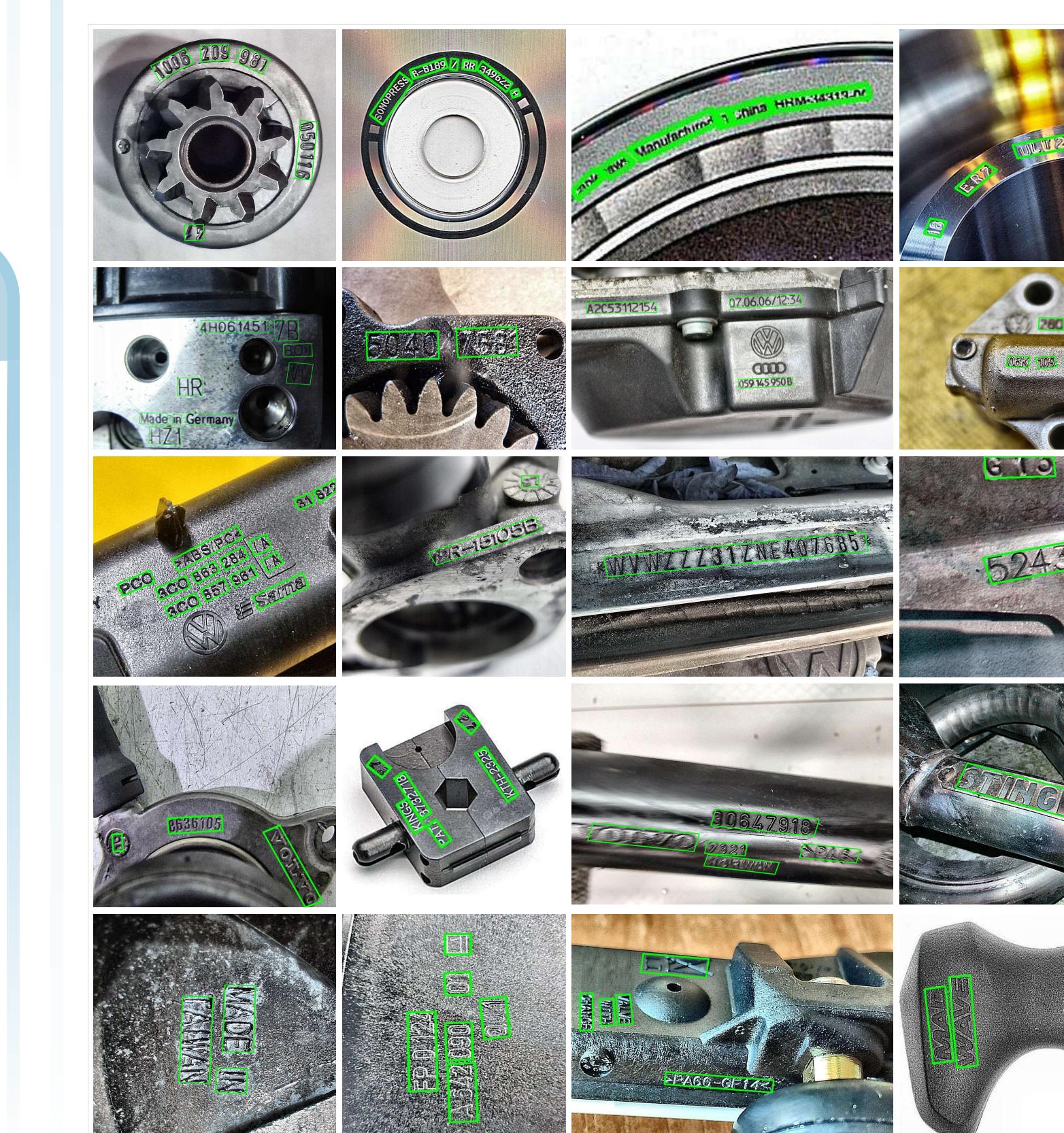


**Analysis of SFF and APR Models:** 1) Scale-sensitive feature fusion. It enhances the feature complementarity and generates robust feature representations with multi-scale texts. 2) Box selection algorithm by multi-scale masks. Background boxes have a low opportunity to be re-corrected, which reduces FP. More foreground boxes (several times than RPN) at each ground-truth box are obtained, which improves the quality of candidates. It promotes positive learning of the refined network and obtains more high-IOU boxes, increasing TP.

**Qualitative Results on the MPSC Dataset and Publicly Available Natural Text Datasets:**

MPSC			MSRA-TD500			ICDAR2013			ICDAR2017-MLT		
Algorithms	Precision (%)	Recall (%)	F-measure (%)	Algorithms	Precision (%)	Recall (%)	F-measure (%)	Algorithms	Precision (%)	Recall (%)	F-measure (%)
EAST	76.33	73.04	74.65	TextSnake†	83.2	73.9	78.3	SegLink*	92	84.4	88.1
Mask R-CNN	85.28	79.25	82.15	PixelLink*	83	73.2	77.8	SSTD	89	86	88
RRPN	81.98	78.91	80.42	RRPN	82	68	74	LOMO	78.8	86	89
PSENet	85.42	78.4	81.76	RRD†	87	73	79	RRD*	92	86	89
PAN	87.07	81.6	84.24	Lyu et al.	87.6	76.2	81.5	PixelLink*	88.6	87.5	88.1
BDN	86.6	77.49	81.79	AS-RPN	84.7	80.4	82.5	RRPN	84	77	80
ContourNet	87.79	81.02	84.27	CRAFT	88.2	78.2	82.9	Melinda et al.	93.9	91.5	92.6
RRPN++	86.73	83.9	85.3	ATTR	85.2	82.1	83.6	FTPN	93.2	91.9	92.5
FCENet	87.13	81.63	84.29	PAN‡	84.4	83.8	84.1	Liu et al.	90.2	86.3	88.2
RFN (ours)	89.3	83.33	86.21	RFN	88.4	80	84	Wei et al.	93.7	87.4	90.4
RFN* (ours)	89.82	84.45	87.05	RFN‡	88.4	87.8	88.1	RFN (ours)	92.5	90.7	91.6

**Qualitative Results on MPSC Dataset:**



**Qualitative Results on Natural Scene Datasets:**

