

# 1. Background & Motivation

The recent studies of Chinese character recognition (CCR) can be roughly divided into two categories: character-based CCR and radical-based CCR.

- Character-based methods can perform well on common Chinese characters with a lot of training data. But they have difficulty in dealing with zero-shot learning problem and perform poorly when handling complex characters
- Radical-based methods has the ability to recognize unseen Chinese characters, However, to recognize more complicated radical structures or learn the composition rules of low-frequency samples.

Therefore, it is important to build a more powerful recognition system for CCR.

- Self-attention mechanism can capture long-range dependencies and the detailed internal pattern
- The Transformer is composed of stacked blocks and aggregates the input context for each block, which naturally provides us with more hierarchical representations

Hierarchical radical structure of an example Chinese character



# 2. The Proposed Method

Architectures of the Transformer-based radical analysis network (RTN) contains two components:

- a dense encoder which takes the image as input to produce a fixed-(1)length context vector;
- (2)a transformer decoder which takes the context vector as input to generate a variable-length symbol sequence.



### > A dense encoder

 $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_L\}, \mathbf{a}_l \in \mathbb{R}^D$ > A transformer decoder  $s_i^n = MA(c_i^{n-1}, C_{<i}^{n-1}, C_{<i}^{n-1})$ 

 $\mathbf{z}_i^n = MA(\mathbf{s}_i^n, \mathbf{A}, \mathbf{A})$ 

 $\mathbf{c}_{i}^{n} = \max\left(0, \mathbf{W}_{1}^{n}\mathbf{z}_{i}^{n} + \mathbf{b}_{1}^{n}\right)\mathbf{W}_{2}^{n} + \mathbf{b}_{1}^{n}$ 

Training Objective

 $\mathbf{P}(\mathbf{Y}|\mathbf{X}) = \prod_{i=1}^{|\mathbf{Y}|} \mathbf{P}\left(\mathbf{y}_{i}|\mathbf{Y}_{< i}, \mathbf{X}\right) = \prod_{i=1}^{|\mathbf{Y}|} \mathbf{P}\left(\mathbf{y}_{i}|\mathbf{c}_{i}^{N}\right)$ 



## 3. Experiments and Results

Our experiments are conducted on both printed Chinese character dataset and natural scene Chinese character dataset.

- $\triangleright$ Experiments on single-font printed Chinese characters
  - Comparison of accuracy rate between RTN and RAN with different caption lengths on different font-style unseen Chinese characters respectively.

| Font Style | RAN(%) |          |       | RTN(%) |          |       |
|------------|--------|----------|-------|--------|----------|-------|
|            | ALL    | $\leq 6$ | > 6   | ALL    | $\leq 6$ | > 6   |
| Song       | 92.21  | 93.65    | 90.78 | 94.54  | 94.93    | 94.16 |
| FangSong   | 91.04  | 91.98    | 90.11 | 94.21  | 94.84    | 93.57 |
| Hei        | 90.41  | 91.34    | 89.50 | 92.79  | 92.41    | 93.11 |
| Kaiti      | 88.57  | 90.59    | 86.58 | 91.31  | 92.96    | 89.67 |

rison of character-level accuracy between RTN and RAN with respect to the frequency of radicals.



- Experiments on natural scene Chinese characters ۶
  - Comparison of the performance of RAN and RTN with the different appearance frequency of character-level categories

| Frequency<br>Categories<br>Samples | $\leq 20$<br>398<br>1128 | $\leq 50 \\ 511 \\ 1229$ | $\leq 100 \\ 335 \\ 1663$ | HF<br>1044<br>48745 | ALL<br>2015<br>52765 |
|------------------------------------|--------------------------|--------------------------|---------------------------|---------------------|----------------------|
| RAN                                | 25.88%                   | 47.92%                   | 65.12%                    | 89.01%              | 85.95%               |
| RTN                                | 41.84%                   | 61.51%                   | 7 <b>1.6</b> 7%           | 89.76%              | 87.51%               |

Comparison of the recognition performance of RTN and RAN with different caption lengths on the CTW valid database

#### Comparison of the recognition performance of RAN and RTN with respect to 6 attributes on the CTW test dataset.

| Model                 | Caption length |          |        |  |
|-----------------------|----------------|----------|--------|--|
|                       | ALL            | $\leq 4$ | > 4    |  |
| RAN                   | 85.95%         | 89.33%   | 82.31% |  |
| RTN                   | 87.51%         | 90.02%   | 84.80% |  |
| Accuracy <sup>↑</sup> | 1.56%          | 0.69%    | 2.49%  |  |

Attributes Training Samples RTN(%) RAN(%) 760107 101393 218560 87.31 73.94 84.57 83.60 78.06 84.25 66.70 85.56 71.55 82.84 71.55 76.17 87.11 63.58 All

192481 199066

65983 6661

|    | Conducion   |
|----|-------------|
| 4. | GOLICIUSION |
|    |             |

✓ We explore the option to improve the capability of RAN by employing the Transformer architecture.

background distorted 3D raised

wordart handwritter

- ✓ The proposed model achieves significant performance improvements on both printed Chinese character database and natural scene Chinese character database.
- ✓ Further analysis proves that RTN is more effective and robust than RAN for recognizing complicated and lowfrequency samples