L DENet: A Holistic Approach to Offline Handwritten Chinese and Japanese Text Line Recognition

ICPR 2020, January 10th-15th, 2021

Huu-Tin Hoang*, Chun-Jen Peng*, Hung Vinh Tran, Hung Le

Cinnamon AI Minato, Tokyo 105-0001, Japan Email: {tin, larry, xing, toni}@cinnamon.is



*Corresponding author

Hoang*, Peng*, Tran, Le, Nguyen



Huy Hoang Nguyen

University of Oulu Oulu, Finland Email: huy.nguyen@oulu.fi



ICPR 2020



2/14

Vast style variety of handwritten texts, especially for Chinese and Japanese





3/14

There are over 47,000 kanji characters! ... and here is the one-hot encoding keyboard.



https://knowyourmeme.com/photos/479661-wtf

Level of decomposition



Difference between our encoding (radicals + basic components) and IDS

- Radicals: stopped at 214 Kangxi radicals and 115 CJK radicals.
- Basic components: If no radical occurred, stopped at leaf nodes



LODEC vs Other encodings





(a) A handwritten sample of the character 號 and the corresponding encoding methods: (b) Cangjie, (c) Ideographic Description Sequences (IDS) and (d) LODEC.

ICPR 2020

LODENet: Overview





Feature Extractor and Sequential

Decoder: CRNN architecture inspired by *Baoguang Shi et.al.**, for predicting radicals

Conversion Network: composed of inception convolutional and sequential layers, for predicting logograms

LODENet overview picture

* Baoguang Shi et.al., An End-to-End Trainable Neural Network for Image-based Sequence Recognition and Its Application to Scene Text Recognition



	CASIA	SCUT-EPT	Private*	
Dataset information				
Language	Chinese	Chinese	Japanese	
Writers	1,080	2,986	-	
Input height h	128	64	64	
Classes	7,355	4,251	2,227	
Amount of data (lines/characters)				
Train	52,229 / 1.4M	40,000 / 1.0M	14,595 / 248K	
Synthesis-Random texts	120,000 / 1.9M	120,000 / 1.9M	-	
Synthesis-Wiki texts	120,000 / 1.9M	120,000 / 1.9M	-	
Test	3,432 / 92K	10,000 / 249K	2,941 / 50K	

Description of 3 datasets used in this study.

*D. Nguyen et al. '19, "Improving long handwritten text line recognition with convolutional multi-way associative memory"

Madal	Character (%)		Radical (%)	
Model	AR ↑	CR ↑	AR ↑	CR↑
CNN + MDLSTM + CTC, Messina et al. '15	73.26	78.30		
CNN + MDirLSTM + CTC, Wu et al. '17	73.65	78.53		
CNN + LSTM + CTC, Y. Zhu et al. '18	75.97	80.26		
CMAM, D. Nguyen et al. '19	74.45	82.14		
LODENet (Ours)	76.61	82.91	77.81	85.40
LODENet + Random texts (Ours)	77.36	84.64	79.27	87.47
LODENet + Wikipedia texts (Ours)	77.61	83.74	79.25	86.33

Testing results on SCUT-EPT.

Model	Character (%)		Radical (%)	
	AR ↑	CR ↑	AR ↑	CR↑
HIT-1, C. Szegedy et al. '15	83.58	86.15		
HIT-2, C. Szegedy et al. '15	86.73	88.76		
Wang et al. '16	88.79	90.67		
CNN + SMDLSTM + CTC, Wu et al. '17	86.64	87.43		
LODENet (Ours)	86.82	87.83	88.51	91.09
LODENet + Random texts (Ours)	92.00	92.89	92.79	95.25
LODENet + Wikipedia texts (Ours)	92.16	93.04	92.77	94.65

Testing results on CASIA.

Madal	Character (%)		
Woder	Validation \downarrow	Test ↓	
CMAM, D. Nguyen et al. '19	17.55	12.99	
LODENet w/o conversion network (Ours)	12.36	6.25	
LODENet (Ours)	11.70	5.47	

CERs on the Japanese dataset.

Ablation studies on different components of our method.

	与其城市建築	有很大相关性
Ground truth	与其城市文化有很大相关性	
CRNN	与其城市菊	(8)
LODENet (Ours)	与其城市文化有有大相关性	(1)
Radical output	与一甘一八土戊了一小1七	1 木月11一大人忄青
	言は若高い	-\$165 物巴
Ground truth	责任在肩头	善假于物也
CRNN	责近"寿高头 (4)	苦物也 (3)
LODENet (Ours)	责任"在肩头 (1)	善假于物也 (0)
	▲□イ―+"+―□□	

LODENet and CRNN on SCUT-EPT test set

12/14

🛞 cınnamon Al

- 1. **LODEC encoding** that can fully represent all logograms and syllabic characters of Chinese and Japanese.
- 2. An **end-to-end training scheme** that can be plugged in any sequential architecture and radical-based encoding method.
- 3. **LODENet architecture** equipped with the conversion network that learns to transcribe Japanese and Chinese contents from radical-based features.
- 4. **SOTA results** on CASIA and SCUT-EPT, and one private Japanese dataset.

Thank You!

Poster Session: PS T4.1 - 1319

Appendix

15/14

Accurate rate (AR) and correct rate (CR), which were used in the ICDAR 2013 competition.

$$CR = \frac{N - Sub - Del}{N}$$
$$AR = \frac{N - Sub - Ins - Del}{N}$$

Character error rate (CER) on the Japanese dataset.

$$CER = \frac{Sub + Ins + Del}{N}$$

Learning Curve

Learning curves of LODENet with different encoding methods on the SCUT-EPT dataset.

18/14

Method	No. of codes	Bijective function	Encoding length Mean & Std.
One-hot	21448	Yes	21448±0
Cangjie	26	No	109±19
IDS	387	No	3618±1722
LODEC (Ours)	520	Yes	1856±691

Table I: The amounts of codes and encoding lengths to represent 21,448 Kanji characters having Unicode correspondences in different encoding methods.

LODENet: Feature Extractor

LODENet: A Holistic Approach to Offline Handwritten Chinese and Japanese Text Line Recognition

19/14

LODENet: Conversion Network

Inspiration: Kanji input method

Problem: There are over 47,000 kanji characters!!

Similar repetitive components appear in different characters

http://202.175.185.186/it-school/mbc/association/sec/it/cj_note.htm

Cangjie input method

- First Chinese input method on QWERTY, invented by Chu Bong-Foo in 1976
- Composed of **24 radicals** associated with **77 auxiliary shapes** (total 101 basic components)

倉頡字母表及輔助字形

Cangjie input method

https://www.wikiwand.com/en/Cangjie_input_method

https://zh.m.	wikipedia.org/ 🔟 i
≡ WIKIPEDI	A C
倉頡輸入法	
ŻД	* * •
ALARRA & NALES TRAIN	PD Ab also be to A Ab about factor about a
倉頡輸入法是一種常 腦之父」美醫的朱邦得 有正體中文版本,原得 腦處理漢字的問題,(儲存、漢字排序等。Я 三軍大學發展中文通 長的蔣緯國為紀念上 1978年將此輸入法面	用的中文輸入法,由有「中文(复先生於1976年創製。初期只 名「形象檢字法」,用以解決備 包括漢字輸入、字形輸出、內面 卡邦復發明此輸入法時正值他 和復發明此輸入法時正值他 和 意語與意讀造字的精神,乃於 新定名為「意顏輸入法」。

ICPR 2020

Ideographic Description Sequences (IDS)

ICPR 2020

- Unicode defines **12 IDCs** (Ideographic Description Characters)
- Different encoding for same component
- Not fully decomposable

Unicode	IDC	Example	IDS
2FF0		相	□□木目
2FF1		志	□士心
2FF2		湘	Ⅲ氵木目
2FF3		糞	□米田共
2FF4		回	
2FF5		回	
2FF6		凶	ШЦХ
2FF7			
2FF8		病	□疒丙
2FF9		戴	□[異
2FFA	The second se	起	□走巳
2FFB		巫	回工从

LODEC for Japanese

25/14

• Codec package for JApanese character DEComposition

- 1. One-hot :輝
- 2. Cangjie : FUBJJ
- 3. IDC_Cangjie : IFUBJJ
- 4. DCs : `'' 🗆 🗆
- 5. IDC_DCs : '' □ □ □

- Hiragana & Katakana decomposition: Isolate dakuten and handakuten

 ホ -> ["ホ"]
 ボ -> ["ホ", "^{*}"]
 ポ -> ["ホ", "[°]"]
- 7. original_DCs
 : `'' □ □ □ □ □

 8. original_IDC_DCs
 : □ `'' □ □ □ □ □

Discussion

- 1. Why do we use 1st Level IDC?
 - Most Kanjis are simple 2-level IDS tree
 - The major space between radicals (described by 1st level IDC) is an important feature to learn

- 1. Difference between our encoding (radicals + basic components) and IDS
 - **Radicals**: We stopped at 214 Kangxi radicals and 115 CJK radicals.
 - **Basic components**: If no radical occured, we stopped at leaf nodes

