

CNN-Based Text Image Super-Resolution Tailored for OCR

Haochen Zhang, Dong Liu, Zhiwei Xiong

CAS Key Laboratory of Technology in Geo-Spatial Information Processing and Application System, University of Science and Technology of China





育創

天寰



Part1 Introduction

Part2 Method

Part3 Experiment

Part4 Conclusion





Part one

Introduction





Note: low-resolution images may hamper the performance of OCR.





PSNR vs OCR accuracy

RESULTS OF DIFFERENT METHODS ON THE ICDAR 2015

IEXISK DATASET				
Method	RMSE	PSNR	MSSIM	OCR (%)
Bicubic	19.04	23.5	0.879	60.64
Lanczos3	16.97	24.65	0.902	64.36
Orange Labs	11.27	28.25	0.953	74.12
Zeyde et al.	13.05	27.21	0.941	69.72
A+	10.03	29.5	0.966	73.1
Synchromedia Lab	62.67	12.66	0.623	65.93
ASRS	12.86	26.98	0.95	71.25
SRCNN-1	7.52	31.75	0.98	77.19
SRCNN-2	7.24	33.19	0.981	76.1
百創				

Two points:

- SR improves the accuracy of OCR.
- The method achieving the highest PSNR does not achieve the highest OCR accuracy.

Contributio

1



We summarize our attempts of improving CNN-based text image SR method to facilitate OCR.

- A new loss function
- Image padding method
- Model combine

Experimental results:







Part two

Method







We use a weighted MSE (WMSE) to <u>emphasize high contrast edges</u> more than others to guide CNN to <u>concern more on the high-frequency image details</u>.

WMSE =
$$\frac{\sum_{i=1}^{m} \sum_{j=1}^{n} \left\| I(i,j) - \hat{I}(i,j) \right\|^2 \times f[\operatorname{grad}(i,j)]}{mn}$$

- I and \hat{I} are the original and super-resolved images, respectively
- *grad* is the gradient magnitude map of the original image.
- $f[\cdot]$ is a certain function to convert gradient magnitude into weight







Model combine



Suppose we have

- Multiple trained networks N_i , i = 1, ..., M
- A set of combination weights w_i , i = 1, ..., M

Combine as
$$\hat{I} = \sum_{i=1}^{M} w_i \times N_i(I_{LR})$$

We do have

育創

天寰

- N_i : Trained by different random initializations
- w_i : Exhaustive search from a finite set



Image padding



In each recursion







Part three

Experiment







Image padding

HR image



No padding



Duplication

Our padding

畏



RESULTS OF DIFFERENT PADDING METHODS

Method	PSNR	MSSIM
No padding	32.18	0.9841
Duplication padding	32.31	0.9842
Our proposed padding	32.42	0.9842



New Loss Function

loss fuction	oss fuction Experiment			DENID (dB)
1055 1001011	Lypenment	accuracy		
	Net 1.1	76.34%	0.9795	31.65
	Net 1.2	76.54%	0.98	31.71
1) f(x)=x^2	Net 1.3	76.41%	0.9796	31.57
	Net 1.4	75.86%	0.9795	31.59
	average	76.29%	0.97965	31.63
	Net 2.1	76.54%	0.9815	31.87
	Net 2.2	76.82%	0.9812	31.80
2) f(x)=x	Net 2.3	75.15%	0.9813	31.91
	Net 2.4	76.44%	0.9811	31.74
	average	76.24%	0.981275	31.83
育創	Net 3.1	77.23%	0.9828	32.00
天寛	Net 3.2	75.15%	0.9821	31.94
嚴 3) f(x)=x^1/2	Net 3.3	76.72%	0.9826	31.94
濟華的	Net 3.4	77.98%	0.9827	31.98
	average	76.77%	0.98255	31.97
题小们	T	able 1 WM9	SE	

Table 2 State-of-art				
Method	RMSE	PSNR	MSSIM	OCR (%)
SRCNN-1	7.52	31.75	0.98	77.19
SRCNN-2	7.24	33.19	0.981	76.1

Table 3 Standard MSE

Method	PSNR	MSSIM
No padding	32.18	0.9841

* The OCR accuracy is measured by using the Tesseract-OCR software.



Model combine

$$\hat{I} = \alpha \cdot N_1(I_{LR}) + (1 - \alpha) \cdot N_2(I_{LR})$$
 $N_1 = Net 3.4$
 $N_2 = Net 3.2$

dataset	method	OCR (%)	OCR of combine(%)	
training	Net 3.2	73.58	74.25	
set	Net 3.4	73.93	74.25	*
testing	Net 3.2	75.15	70 10	$\alpha = 0.67$
set	Net 3.4	77.98	/8.10	↓







Part four

Conclusion





育

創

寰



- A new loss function to guide the CNN to focus on high-First frequency image details.
- Image padding method to refine the image boundaries Second during CNN based SR.







Thank you for listening!

育 創 天寰 下字





Why the results of using VDSR on the text images are not satisfactory, especially at image boundaries ?

- Text images in the TextSR dataset are much smaller compared with natural images.
- VDSR is very deep (20 layers), the receptive field of the last convolutional layer is 41×41 .





CDAR 2015 TextSR dataset



The training set consists of 567 pairs of HR-LR grayscale images and the ground-truth OCR results

The test set consists of 141 pairs of HR-LR images as well as their ground-truth OCR results

OCR results include English letters, numbers, and 14 special characters such as ",".



Implementation Details





10 layer vs 20 layer

Table 1 10 layer VDSR

Learning rate	bicubic	VDSR
10^-4	23.5011	30.0127
10^-5	23.5011	29.9136
10^-6	23.5011	29.8435
10^-4>10^-5	23.5011	30.0622
10^-5>10^-6	23.5011	29.9276
10^-4>10^-7	23.5011	30.0061
10^-5>10^-8	23.5011	29.8794

育

創

宇

天寰

Table 2 20 layer VDSR			
Method	PSNR	MSSIM	
No padding	32.18	0.9841	

SR vs original HR



Après utilisation du shampooing force & densité, pellicules visibles, usage régulier,

Label: Apres_utilisation_du_shampooing_force_&_densite,_pellicules_visibles,_usage_regulier
SR: Apms_utilisation_du_shampooing_i_&_densite,_pellicules_visibies,_usage_rigulier,
HR: iAprths_utilisation_du_shampooing_IQ_&_densite_peilicules_visibles,_usage_r&gulier.

