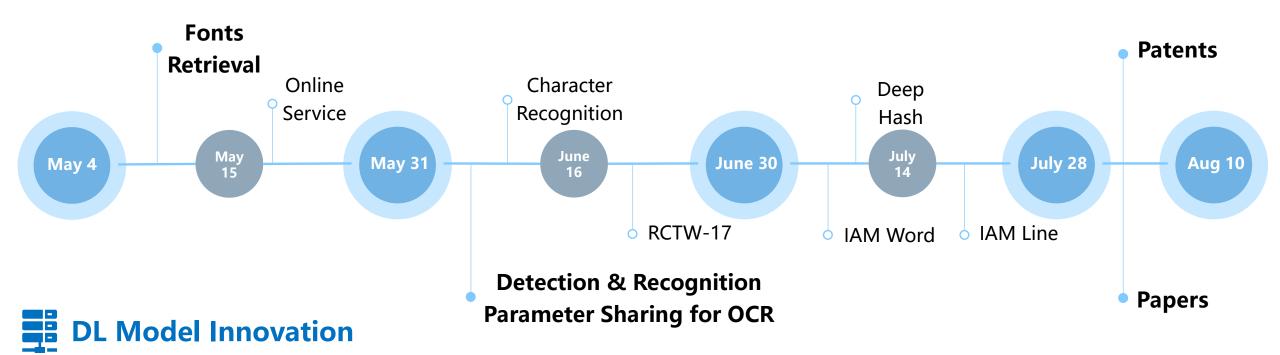


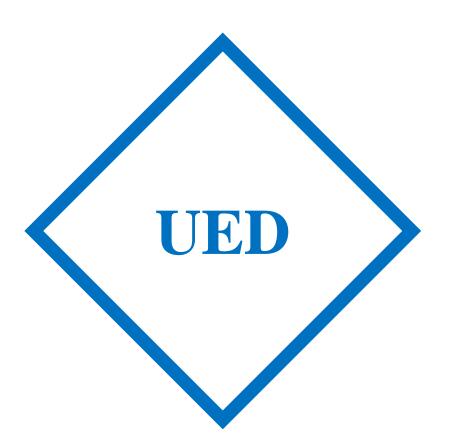
实习总结汇报

Bin Xu-Large Scale Learning-Deep Learning



UED Intelligent Design

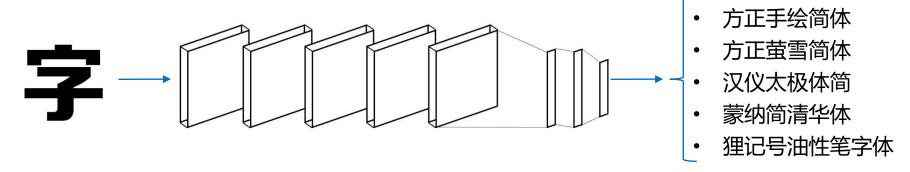






Fonts Retrieval

Model: we use the CNNs for the retrieval task.



• Dataset: we build two datasets according to the font files provided by International UED team.

| Font Classes | Training Set | Testing Set | Validation set |
|--------------|--------------|-------------|----------------|
| 506 | 20000k | 20K | 20K |
| 1411 | 50000k | 20K | 20K |

Examples:













Training Model — we train and evaluate three classification models: ResNet, GoogleNet and VGG on 506 and 1411 font classes dataset, respectively.

• 506 font classes:

| Model | val-top5 | val-top1 | test-top5 | test-top1 | Iteration | |
|--------------|----------|----------|-----------|-----------|-----------|--|
| ResNet-18 | 0.96985 | 0.5235 | 0.96645 | 0.49785 | 210000 | |
| GoogleNet-4c | 0.96015 | 0.4886 | 0.9578 | 0.4735 | 180000 | |
| VGG-16 | 0.97835 | 0.5161 | 0.97645 | 0.51085 | 185000 | |

• 1411 font classes:

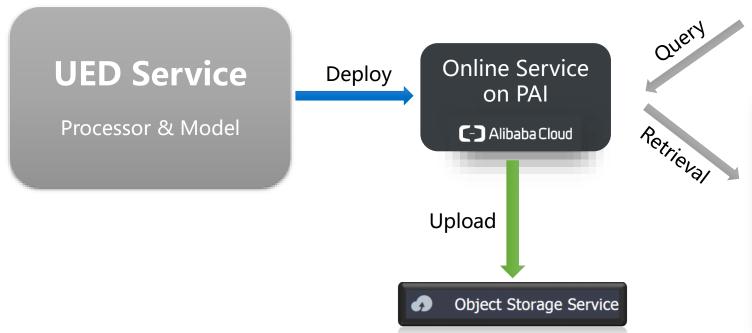
| Model | test-top20 | test-top10 | test-top5 | test-top1 | Iteration |
|--------------|------------|------------|-----------|-----------|-----------|
| ResNet-18 | 0.98685 | 0.97055 | 0.90355 | 0.4923 | 205000 |
| GoogleNet-4c | 0.9873 | 0.9697 | 0.90205 | 0.48925 | 370000 |
| VGG-16 | 0.98405 | 0.96275 | 0.88445 | 0.461 | 280000 |





Fonts Retrieval Service

• We deploy the online service on **PAI**.





WiKi: http://gitlab.alibaba-inc.com/jason.xb/online_processor/wikis/home





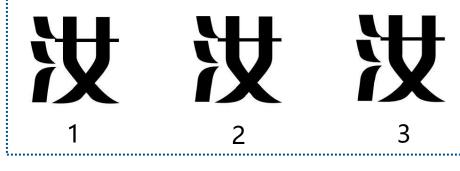
Classification or Retrieval?

• The classification models output a multiclass categorical probability distribution by the softmax function, and are trained by the Cross Entropy loss function.

$$y_{c} = \varsigma(z)_{c} = \frac{e^{z_{c}}}{\sum_{i=1}^{C} e^{z_{d}}} \qquad \text{for } c = 1...C \qquad \qquad \xi(T, Y) = \sum_{i=1}^{n} \xi(t_{i}, y_{i}) = -\sum_{i=1}^{n} \sum_{i=c}^{C} t_{ic} \cdot \log(y_{ic})$$

 We notice that the classification models learn the nonlinear functions from each font class separately, but ignore the inner connection between two similar fonts.



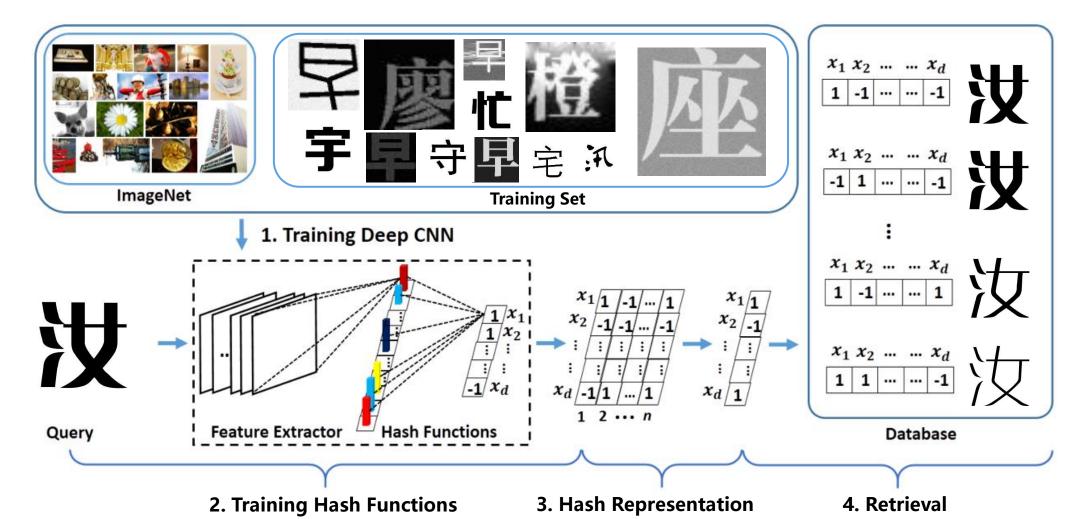








Deep Hash Method — we employ hash coding method in our retrieval model.

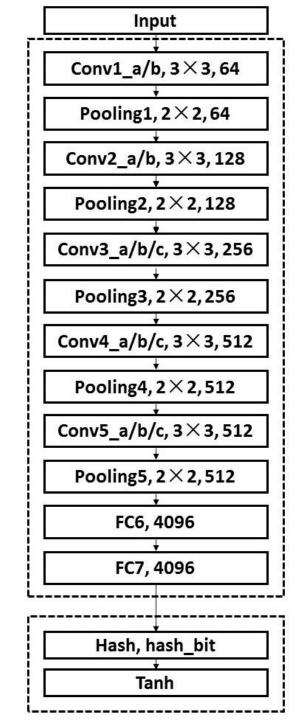






Proposed Network

- We use the released VGG (Visual Geometry Group) model which is trained on ImageNet dataset for convenience.
- Previous studies show that the 4096-dimensional features of Layer7 perform better than many handcrafted features. In our network, the Layers1-7 are used as the Feature Extractor Network.
- We set Hash (full connected) layer and Tanh activation layer behind as the Hash Network.





Deep Hash Representation

• Hash code of deep features:

$$\mathbf{p}_i = h_f(\mathbf{f}_i),$$
 $h_f : \mathbf{F} \to (-1, +1)^S$ $\mathbf{c}_i = h(\mathbf{f}_i) = sgn(\mathbf{p}_i),$ $h : \mathbf{F} \to \{-1, +1\}^S$

Hamming distance between two codes:

$$d\left(\boldsymbol{c}_{i},\boldsymbol{c}_{j}\right) = \frac{S - \boldsymbol{c}_{i}^{T} \boldsymbol{c}_{j}}{2} \longrightarrow$$

Hash Representation of Font:

$$\bar{c} = \frac{1}{N} \sum_{i=1}^{N} c_i$$

$$d(\mathbf{c}_{i}, \mathbf{c}_{j}) = \frac{S - \mathbf{c}_{i}^{T} \mathbf{c}_{j}}{2} \longrightarrow d(\mathbf{c}_{1}, \mathbf{c}_{2}) = \frac{3 - (1 * 1 + 1 * (-1) + (-1) * 1)}{2}$$

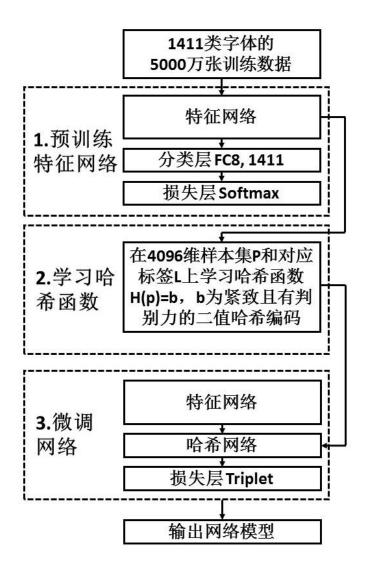
$$= \frac{3 - (-1)}{2} = 2$$





Training Procedure

- The training procedure of our network contains three steps:
 - pre-train the Feature Extractor Network
 - train the Hash Functions
 - fine-tune the Unified Network

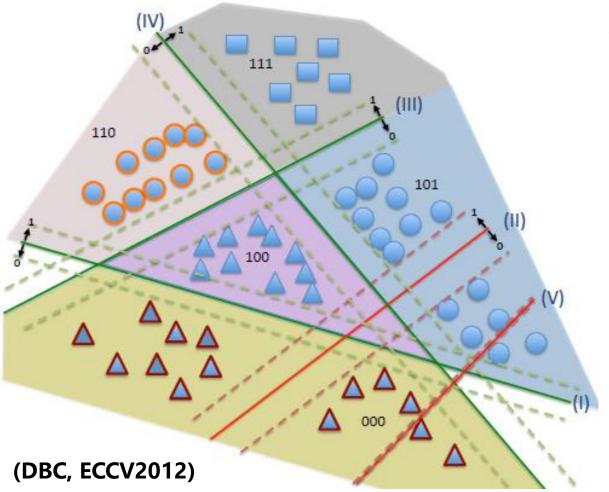




• Deep features:
$$P = [p_1, p_2, ..., p_N] \in \square^{T \times N}$$

• Hash functions:
$$W = [w_1, w_2, ..., w_S] \in \square^{T \times S}$$

• Hash coding:
$$\mathbf{B} = sgn(\mathbf{W}^T \mathbf{P}) = [\mathbf{b}_1, \mathbf{b}_2, ..., \mathbf{b}_N] \in \{+1, -1\}^{S \times N}$$



• To learn hash functions, **discrimination** constraint and **stability** constraint are taken into account.

Objective:
$$\min_{W, B} \frac{1}{2} ||W||^2 + \mu \sum_{i=1}^{S} \sum_{j=1}^{N} \xi_{ij} + g(B)$$

s.t.
$$\boldsymbol{B}_{ij}(\boldsymbol{w}_{i}^{T}\boldsymbol{p}_{j}) \geq 1 - \xi_{ij},$$

 $\xi_{ij} \geq 0,$
 $\boldsymbol{B} = sgn(\boldsymbol{W}^{T}\boldsymbol{P})$

$$g(\boldsymbol{B}) = \sum_{\substack{p=1,\\\boldsymbol{b}_i,\boldsymbol{b}_j \in \boldsymbol{B}_p}}^{C} d(\boldsymbol{b}_i,\boldsymbol{b}_j) - \lambda \sum_{\substack{p=1,\\\boldsymbol{b}_i \in \boldsymbol{B}_p}}^{C} \sum_{\substack{q=1,q \neq p,\\\boldsymbol{b}_j \in \boldsymbol{B}_q}}^{C} d(\boldsymbol{b}_i,\boldsymbol{b}_j),$$

Fine-tune

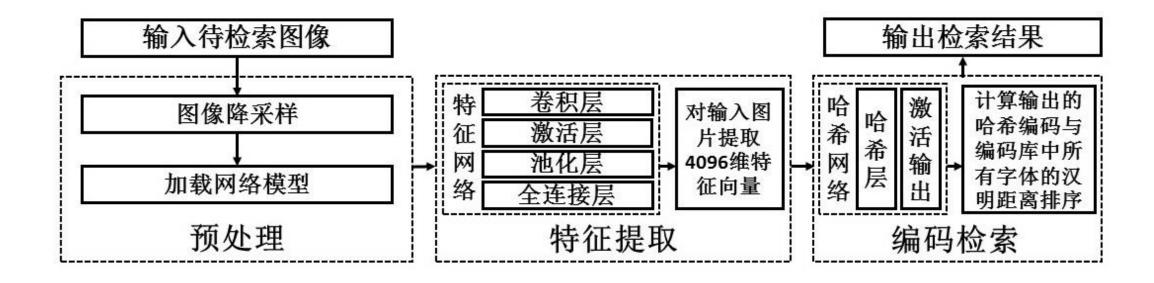
- We use the **Triplet Ranking Loss** similar to *Hamming Distance Metric Learning* (NIPS 2012), to fine-tune our network.
 - sample p is more similar to p_+ than to $p_ \longrightarrow$ $d(p, p_+) + \zeta < d(p, p_-)$
 - loss function $L(\boldsymbol{p}, \boldsymbol{p}_{+}, \boldsymbol{p}_{-}) = \max(d(\boldsymbol{p}, \boldsymbol{p}_{+}) d(\boldsymbol{p}, \boldsymbol{p}_{-}) + \zeta, 0)$

• objective function $\min_{\boldsymbol{W}_h, \boldsymbol{W}} \sum_{i=1}^{C} \sum_{\substack{\boldsymbol{p}, \boldsymbol{p}_+ \in \boldsymbol{F}_i, \ \boldsymbol{p}_- \in \boldsymbol{F}_j \\ \boldsymbol{p} \neq \boldsymbol{p}_+}} \sum_{\substack{\boldsymbol{p}_- \in \boldsymbol{F}_j \\ j \neq i,}} \max \left(d\left(\boldsymbol{p}, \boldsymbol{p}_+\right) - d\left(\boldsymbol{p}, \boldsymbol{p}_-\right) + \zeta, 0 \right)$





Patent: 《一种基于深度卷积神经网络与哈希编码的字体检索方法和装置》







• We train and evaluate our retrieval model on 506 font classes dataset:

| Model | test-top20 | test-top10 | test-top5 | test-top1 | Iteration | |
|--------------|------------|------------|-----------|-----------|-----------|--|
| ResNet-18 | - | - | 0.92655 | 0.43925 | 105000 | |
| GoogleNet_4c | - | - | 0.92235 | 0.43945 | 195000 | |
| VGG_16 | - | - | 0.92645 | 0.43085 | 185000 | |
| VGG+Hash | 0.977 | 0.965 | 0.929 | 0.474 | 280000 | |

• Retrieval results:



Query

拨拨拨汝汝汝汝汝

......





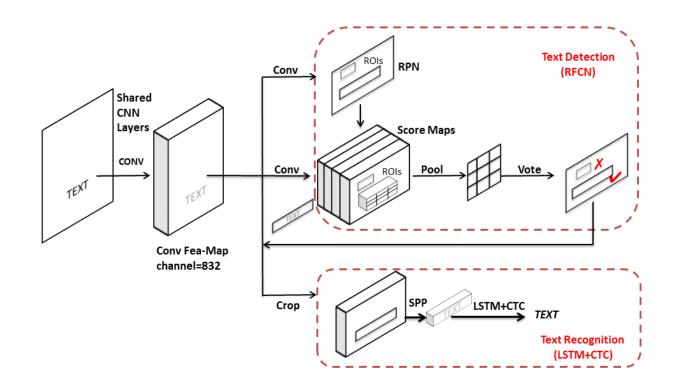
DL Model Innovation



Detection & Recognition Parameter Sharing for OCR — share the parameter for detection and that of recognition so that the time for recognition can be reduced

Two types of training approach:

- Separate training: detector and recognizer trained separately, use the same parameters for shared layers.
- End2end training: detection loss and recognition loss get summed as the whole loss to update the detector and recognizer, respectively.





Parameter Sharing – evaluation on public datasets

ICDAR2017 RCTW-17



Introduction

RCTW-17 is a competition on reading Chinese Text in images. For training and testing, we provide a large-scale dataset that consists of various kinds of images, including street views, posters, menus, indoor scenes, and screenshots. See Dataset and Tasks for details.



- The dataset is built by Huazhong University of Science and Technology.
- 12000 images: 8,000 training, 4,000 testing





 We achieve comparable detection results on testing image (2,000 images divided from RCTW-17 train set).

| F-measure Rank | Institute | Dataset | F-meausre | Precision | Recall |
|-------------------|----------------------------------|--------------------------------|-----------|-----------------------|------------------------|
| 1 | Peking University | | 0.6610 | 74.38% | 59.48% |
| 6 | Shanghai Jiao Tong University | Test set | 0.5760 | 74.06% | 47.12% |
| 9 | Sogou Inc. | (4,229 images) | 0.4961 | 65.43% | 39.95% |
| # | Organizing team (CVPR 2017) | | 0.5278 | 76.03% | 40.43% |
| # | Ours | Train subset (2,000 images) | 0.5421 | 75.65% (7219/9542) | 42.24% (7219/17089) |

 We also test our model on RCTW-17 Test set (annotations not available).



Ground truth Our Train subset(2,000 images)



Our detection results

Our detection results

Test set(4,229 images)



RCTW-17

• Our recognition model performs well on most line samples, but can not recognize characters in declining or vertical samples.





Good case Bad case

The detection and recognition model need to be improved, not finished yet.



IAM (University of Bern)

• The IAM Handwriting Database contains forms of handwritten English text which can be used to train and test handwritten text recognizers and to perform writer identification and verification experiments.

Characteristics

The IAM Handwriting Database 3.0 is structured as follows:

- 657 writers contributed samples of their handwriting
- 1'539 pages of scanned text
- 5'685 isolated and labeled sentences
- 13'353 isolated and labeled text lines
- 115'320 isolated and labeled words

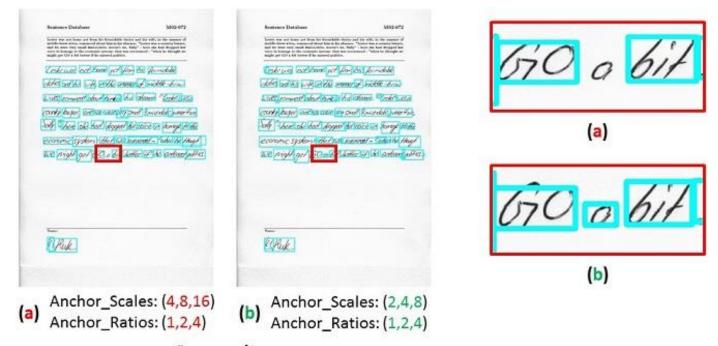
Large Writer Independent Text Line Recognition Task
 This task consists of a total number of 9'862 text lines.

| Set Name | Number of Text Lines | Number of Writers |
|--------------|----------------------|-------------------|
| Train | 6'161 | 283 |
| Validation 1 | 900 | 46 |
| Validation 2 | 940 | 43 |
| Test | 1'861 | 128 |
| Total | 9'862 | 500 |



IAM — Word Level

 We adjust the anchor parameters of the Region Proposal Network (RPN) in R-FCN model for small word samples.



Our results

• We notice that there are a lot of punctuations in IAM, which are easily missed by the detection model, so the recall does not improve much.

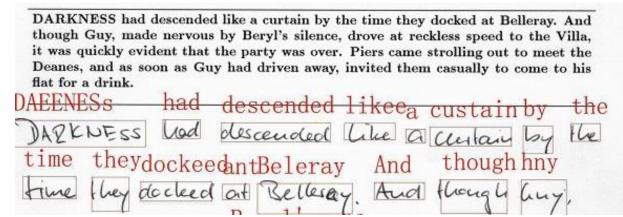


IAM — Word Level

 We train our LSTM+CTC recognition model on the ground-truth bounding boxes, and obtain comparable results.

| Model | Test (17,9 | 83 words) | Validation (19,392 words) | | |
|------------------------|------------|-----------|---------------------------|-----------|--|
| | Word_acc% | Char_acc% | Word_acc% | Char_acc% | |
| MDLSTM-RNNs, NIPS 2016 | 75.4 | 92.1 | 82.3 | 95.1 | |
| Ours (Iter 106,000) | 71.38 | 88.75 | 75.75 | 87.82 | |

• Due to the punctuation and biased detection results, our end-to-end word accuracy can only reach 51.84%.





IAM — Line Level

• In order to recover the missed punctuations, we train our model on line level samples. We also adjust the anchor parameters and obtain accurate detection results.

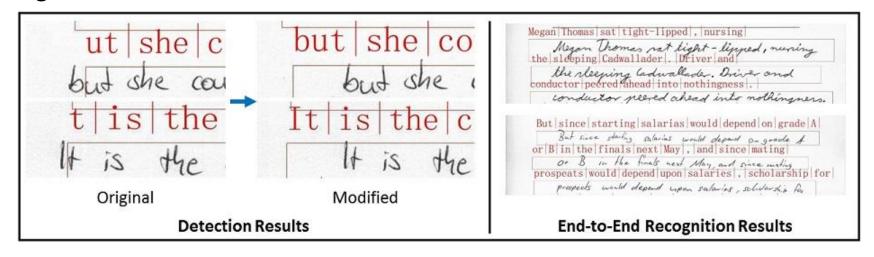
| Threshold | ANCHOR RATIO | ANCHOR SCALES | Recall | Precision | F-Score | False |
|-----------|---------------|---------------|--------|-----------|---------|----------|
| THESHOLD | AIVENON_NATIO | ANCHON_SCALES | recuii | | | Negative |
| IOU>0.5 | (2, 4, 8) | (8, 16, 32) | 0.9881 | 0.9663 | 0.9771 | 22 |
| Score>0.5 | (4, 8, 16) | (4, 8, 16) | 0.9919 | 0.9887 | 0.9903 | 15 |
| 30016>0.3 | (4, 8, 16) | (8, 16, 32) | 0.9940 | 0.9871 | 0.9906 | 11 |
| IOU>0.5 | (4 9 16) | (9 16 22) | 0.9967 | 0.9411 | 0.9681 | 6 |
| Score>0.1 | (4, 8, 16) | (8, 16, 32) | 0.5507 | 0.5411 | 0.9661 | 0 |

• As to recognition model, we adjust the pooling size in the Conv3 and Conv4 to reduce of the loss of information during pooling ROI features.



UED — Line Level

• When analyzing the errors, we notice that the original detection results might miss characters on the boundary of line samples, so we modify the detection results before inputting them into LSTM.



- Our end-to-end recognition word accuracy can reach 76.27% (separate training), and our model can save over 30% of computation time with single thread.
- We are conducting end-to-end model training on the IAM dataset, and we are preparing the paper.

WiKi: https://gitlab.Alibaba-inc.com/sensi.zg/Detection Recognition Parameter Share/wikis/VerificationResult

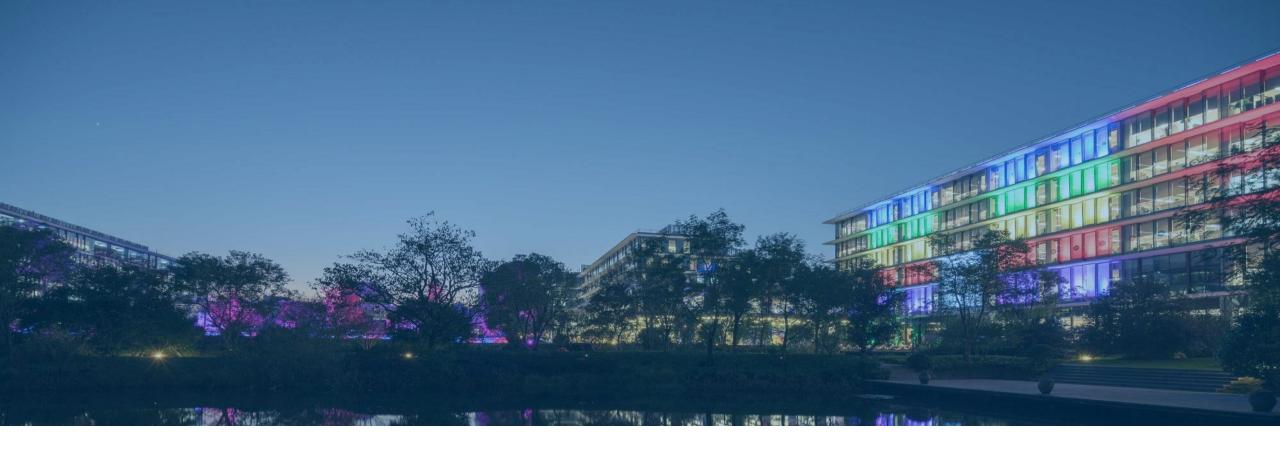




- Fonts Classification CNN Models
- Fonts Retrieval Online Service
- Deep Hash Retrieval Model
- Patent and Paper of Fonts Retrieval Model

DL Model Innovation

- Evaluation Experiments on Public Datasets: RCTW-17、IAM
- Comparable End-to-End Detection & Recognition Results
- Paper of Detection & Recognition Parameter Sharing Model



Thank You

2017.8.10