

:TP391

10406_081203_2004081200109_LW



XXXXXX

:TP391

10001
2004081200109

()

XXXXXXXXXXXXXXXXXXXX

:

:

:

:

:

:

:

Research
Based on Transformer and CNN

A Dissertation

Submitted for the Degree of Master

On the Computer Science and Technology

By Long Liu

Under the Supervision of

A.Prof. Cihui Yang

School of Information Engineering

Nanchang Hangkong University, Nanchang, China

May, 2026

XX
XX
XX
XXX XXX XXX

Abstract

XX
XX
XX

Keywords:XXX XXX XXX

.....	I
Abstract	II
.....	III
1	1
1.1	1
1.2	1
2 XXXXXXXXXXXX.....	1
2.1 XXXXXXXXXXXX.....	1
.....	2
.....	4
.....	5

1

1.1

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

1.2

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

2 XXXXXXXXXXXX

2.1 XXXXXXXXXXXX

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
XXXXXXXXXXXXXXXXXXXXXXXXXXXX^[1] ^[2] ^[3] ^[4]
^[5] ^[6] ^[7] ^[8] ^[9–11] ^[12,13] ^[14]

-
-
- [1] BODLA N, SINGH B, CHELLAPPA R, et al. Soft-nms —improving object detection with one line of code[J/OL]. 2017 IEEE International Conference on Computer Vision (ICCV), 2017: 5562-5570. <https://api.semanticscholar.org/CorpusID:15155826>.
- [2] DALAL N, TRIGGS B. Histograms of oriented gradients for human detection [C/OL]//2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05): Vol. 1. San Diego, CA, USA: IEEE, 2005: 886-893. DOI: [10.1109/CVPR.2005.177](https://doi.org/10.1109/CVPR.2005.177).
- [3] ZEILER M D, FERGUS R. Visualizing and understanding convolutional networks [C]//FLEET D, PAJDLA T, SCHIELE B, et al. Computer Vision – ECCV 2014. Cham: Springer International Publishing, 2014: 818-833.
- [4] WANG X, ZHANG R, KONG T, et al. Solov2: Dynamic, faster and stronger: abs/2003.10152[A/OL]. 2020. <https://api.semanticscholar.org/CorpusID:214611772>.
- [5] WANG C Y, BOCHKOVSKIY A, LIAO H Y M. Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors[C/OL]//2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). 2023: 7464-7475. DOI: [10.1109/CVPR52729.2023.00721](https://doi.org/10.1109/CVPR52729.2023.00721).
- [6] WANG C Y, MARK LIAO H Y, WU Y H, et al. Cspnet: A new backbone that can enhance learning capability of cnn[C/OL]//2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). 2020: 1571-1580. DOI: [10.1109/CVPRW50498.2020.00203](https://doi.org/10.1109/CVPRW50498.2020.00203).
- [7] LIU Y, JIA Q, FAN X, et al. Cross-SRN: Structure-preserving super-resolution network with cross convolution[J/OL]. IEEE Transactions on Circuits and Systems for Video Technology, 2022, 32(8): 4927-4939. DOI: [10.1109/TCSVT.2021.3138431](https://doi.org/10.1109/TCSVT.2021.3138431).
- [8] BRIDLE J S. Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition[C]//SOULIÉ F F, HÉRAULT J. Neurocomputing. Berlin, Heidelberg: Springer Berlin Heidelberg, 1990: 227-236.
- [9] ULYANOV D, VEDALDI A, LEMPITSKY V S. Instance normalization: The missing ingredient for fast stylization: abs/1607.08022[A/OL]. 2016. <https://api.semanticscholar.org/CorpusID:16516553>.
- [10] KASEM M M, ABDALLAH A, BERENDEYEV A, et al. Deep learning for table detection and structure recognition: A survey: abs/2211.08469[A/OL]. 2022. <https://api.semanticscholar.org/CorpusID:253553399>.
- [11] ZHONG X, SHAFIEIBAVANI E, JIMENO YEPES A. Image-based table recognition: Data, model, and evaluation[C]//VEDALDI A, BISCHOF H, BROX T, et al.

Computer Vision – ECCV 2020. Cham: Springer International Publishing, 2020: 564-580.

- [12] MA C, SUN L, ZHONG Z, et al. Relatext: Exploiting visual relationships for arbitrary-shaped scene text detection with graph convolutional networks [J/OL]. Pattern Recognition, 2021, 111: 107684. <https://www.sciencedirect.com/science/article/pii/S0031320320304878>. DOI: <https://doi.org/10.1016/j.patcog.2020.107684>.
- [13] BOLYA D, ZHOU C, XIAO F, et al. Yolact: Real-time instance segmentation [C/OL]//2019 IEEE/CVF International Conference on Computer Vision (ICCV). 2019: 9156-9165. DOI: [10.1109/ICCV.2019.00925](https://doi.org/10.1109/ICCV.2019.00925).
- [14] KHANAM R, HUSSAIN M. Yolov11: An overview of the key architectural enhancements: abs/2410.17725[A/OL]. 2024. <https://api.semanticscholar.org/CorpusID:273532028>.

-
- :
- 1.
 2. AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA

- :
1. XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX-
YYY
 2. BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB

- :
1. XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
YYYYYY
 2. CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
