

GE ZHANG

5 Yiheyuan Road, Haidian District, Beijing, China 100871
(+86)15957610055 ◊ zhangge9194@pku.edu.com

EDUCATION

Peking University, Beijing, China

September 2016 - Present

B.S., Department of Intelligence Science, School of Electronics Engineering and Computer Science
GPA: 3.2/4.0

PUBLICATIONS

1. Zijun Sun*, **Ge Zhang***, Junxu Lu, and Jiwei Li (*: equal contribution)
LOP-OCR: A Language-Oriented Pipeline for Large-chunk Text OCR
2019 ([preprint](#))
2. Anonymous Author(s)
CORAL: COde RepresentATion Learning with Weakly-Supervised Transformers for Analyzing Data Analysis
under submission to SIGKDD 2020([Blog](#))

EXPERIENCE

BData Lab, University of Washington, Seattle

June 2019-Present

Research Intern (visiting student)

- Progress towards a paper on Notebook Analysis, intended for publication at **KDD 2020**.
- Supervisor: Dr. **Tim Althoff**, Assistant Professor, University of Washington, Seattle
- Collaborators: **Mike Merrill**, **Yang Liu**, **Jeffrey Heer**
- Worked on multiverse analysis of public Jupyter Notebooks to help people better understand data science process by detecting decision points and giving alternatives.
- Created a dataset of decision points and alternatives by clustering based on similarity of functions(including signals of return values, arguments, cooccurrence and sommon subtokens).
- Proposed a novel weakly supervised transformer architecture for computing joint representations of data science code from both abstract syntax trees and natural language annotations.
- Presented a new classification task for labeling computational notebook cells as stages in the data analysis process (i.e., data import, wrangling, exploration, modeling, and evaluation).
- Annotated cells with one of the above stages for 100 data science Jupyter notebooks and produced a standardized rubric for qualitative coding.

Shannon.ai

Nov 2018-May 2019

Full-time research intern

- Supervisor: Dr. **Jiwei Li**, Chief Executive Officer, Shannon.ai
- Collaborators: Junxu Lu, Zijun Sun
- Designed a model based on seq2seq with auxiliary image information to help scene text recognition tasks.
- Extended the model to a system to improve the performance of OCR, increasing the accuracy significantly from 77.9 to 88.9.
- (As a product of the company) Extracted structured data from prospectus(PDFs) and helped design a data structure to format long text used for information extraction. Reached 0.95 accuracy on the whole corpus.

- Supervisor: Dr. **Xiaojun Wan**, Professor, Peking University
- Crawled data to construct a fiction-story summarization dataset.
- Used SVM with designed features(length, sentence representation) to summarize public online fiction-stories.

AWARDS

2nd Prize of ACM Competition, Peking University, **2017**

People's Choice Prize, Google Girls' Hackathon (as team leader), **2019**

SELECTED PROJECTS

Multiverse Analyses on Jupyter Notebook

Working towards **KDD 2020**

With the neural generative model backing up, the demo allowed people to upload jupyter notebook, highlighted decision points and suggested alternatives. Allowed people to dynamically add comments and feedback.

Helped data scientists better understand data analyses and provided a platform to better define decision points in data science process.

URL: <http://darkwing.cs.washington.edu:5000> (Closed Beta)

Splendor AI (Google Girls' Hackathon)

An AI bot of Splendor, which could beat most human players and won People's Choice Prize in the hackathon. This AI bot computes scores based on weighted custom features, which also had a try on reinforcement learning.

Github: <https://github.com/AshleyZG/splendor>

LOP-OCR: A Language-Oriented Pipeline for Large-chunk Text OCR

Together with Jiwei Li, Zijun Sun, Junxu Lu.

LOP-OCR, a language-oriented pipeline to help scene text recognition tasks based on seq2seq models with auxiliary image information.

Improve the performance of the CRNN-based OCR models, increasing sentence-level accuracy from 77.9 to 88.9 and position-level accuracy from 91.8 to 96.5.

Love and Hate

A demo of a combinatorial game named Love-and-Hate. In this demo, we generated an optimal way to play the game perfectly in various situations. Also designed with a user friendly GUI.

Github: <https://github.com/AshleyZG/Love-and-Hate>

COURSES

Mathematical Analysis (89), Advanced Algebra (85), Introduction to Computing (85), Computational Game Theory (UC Berkeley) (97), Information Theory (86), Natural Language Processing with Deep Learning (Stanford, mooc)

SKILLS

Chinese(native), English (fluent)

Python, c/c++, Pytorch, Javascript(D3 library)

Familiar with Linux system.