

# Dongjie Cheng

E-Mail: [chengdongjiedev@outlook.com](mailto:chengdongjiedev@outlook.com)  
Homepage: [Dongjie Cheng \(dongjie-cheng.github.io\)](https://dongjie-cheng.github.io)

---

## Education

<b>SICHUAN UNIVERSITY</b> <i>Artificial Intelligence</i>	Chengdu, China Undergraduate, Senior Year
<b>GPA (Compulsory/Overall)</b>	<b>3.88/4, 3.78/4</b>
<b>Rank</b>	<b>4/48</b>
<b>CET-6</b>	<b>636</b>

---

## Honors

<b>National Scholarship (1/48 that year)</b>	2023.11
<b>National Third Prize, “China Software Cup”</b>	2023.08
<b>Regional Second Prize, “RoboMaster North Region Competition”</b>	2023.06
<b>Provincial Second Prize, “National College Mathematics Competition”</b>	2021.12
<b>School Outstanding Student</b>	2022.11
<b>School First-Class Scholarship</b>	2023.11

---

## Research Publications

<i>TV-SAM: Increasing Zero-Shot Segmentation Performance on Multimodal Medical Images Using GPT-4 Generated Descriptive Prompts Without Human Annotation</i>	<b>Co-First Author, Accepted</b> , Big Data Mining and Analytics (JCR Q1, IF=7.7) <a href="#">ArXiv, abs/2402.15759</a>
<i>Calibrated Self-Rewarding Vision Language Models</i>	<b>Co-First Author</b> , Submitting to NeurIPS-2024 Main Track (the short version is presented in ICML 2024 FM-Wild Workshop) <a href="#">ArXiv, abs/2405.14622</a>
<i>Evaluating Hallucination in Text-to-Image Diffusion Models with Scene-Graph based Question-Answering Agent</i>	<b>Co-First Author</b> , Submitting to NeurIPS-2024 D&B Track
<i>SAM on Medical Images: A Comprehensive Study on Three Prompt Modes</i>	<b>Co-First Author, Cited by:78</b> <a href="#">ArXiv, abs/2305.00035</a>

---

## Experience

<b>WEST CHINA HOSPITAL – BIG DATA CENTER</b> <i>Research Assitant of Dr.Kang Li’s Lab</i>	Chengdu, Sichuan February, 2023 – March, 2024
<b>UNC-CHAPEL HILL</b> <i>Remote Intern of Dr.Huaxiu Yao’s Lab</i>	Remote March, 2024 – Present

<b>VLM project</b>	Key Project Member
--------------------	--------------------

Our work addresses these challenges by proposing the Calibrated Self-Rewarding (CSR) approach, which enables the model to self-improve by iteratively generating candidate responses, evaluating the reward for each response, and curating preference data for fine-tuning. In the reward modeling, we employ a step-wise strategy and incorporate visual constraints into the self-rewarding process to place greater emphasis on visual input. Empirical results demonstrate that CSR enhances performance and reduces hallucinations across ten benchmarks and tasks, achieving substantial improvements over existing methods by 7.62%. Our empirical results are further supported by rigorous theoretical analysis, under mild assumptions, verifying the effectiveness of introducing visual constraints into the self-rewarding paradigm.

I was responsible for the specific implementation and optimization of the CSR method, as well as core tasks such as DPO training and SFT training for VLM.

- *Calibrated Self-Rewarding Vision Language Models*

**SAM project**

Key Project Member

In the SAM project, We proposed using large models to generate descriptions for segmentation targets, feeding theses descriptions to the detection model to produce bounding boxes for SAM, thereby achieving zero-shot segmentation.

I was responsible for conceiving and implementing specific experiments. Firstly, I completed the evaluation of the SAM model on multiple modalities medical datasets. Then I verified the effectiveness of the improvement method driven by LLM (Large Language Models).

The results show that the improved method performs well under zero-shot conditions, outperforming the GSAM (GLIP+SAM) baseline on most datasets. The project ultimately resulted in two papers, of which I am a co-first author.

- SAM on Medical Images: A Comprehensive Study on Three Prompt Modes.*
- TV-SAM: Increasing Zero-Shot Segmentation Performance on Multimodal Medical Images Using GPT-4 Generated Descriptive Prompts Without Human Annotation*

**T2i-Eval project**

Key Project Member

In the T2i-Eval project, we proposed a method combining Scene Graph and Graph QA to score the quality of generated images, conducting a comprehensive evaluation of images from perspectives such as object omission, attribute inaccuracies, relational errors, and hallucinations.

I was responsible for generating evaluation dataset images, the specific design and experimentation of the Scene Graph part, achieving the construction of Scene Graphs through the use of GroundingDINO+BLIP VQA.

We constructed a human-evaluated dataset containing 12,000 images from 1,000 prompts and validated the effectiveness of our method. Compared with human evaluations, our Pearson and Kendall correlation coefficients surpassed those of T2ICompbench(Neurips 2023). This project ultimately resulted in one paper, for which I am a co-first author.

- Evaluating Hallucination in Text-to-Image Diffusion Models with Scene-Graph based Question-Answering Agent*

Representative Courses	
Math Courses:	
Matrix Analysis	100
Theory of Optimization (Convex Optimization)	94
Discrete Mathematics	93
Programming Courses:	
Programming Language-I (C)	95
Programming Language-II (C++)	99
Python Programming for Artificial Intelligence	96
Data Structure and Algorithm Analysis	94
Robotics Programming with ROS	98

- Skills
- Proficient in programming languages including Python, C++, and C.
  - Experienced with deep learning libraries such as Hugging Face and PyTorch.
  - Skilled in version control with Git and proficient in using the Linux shell.
  - Equipped with basic theoretical knowledge about machine learning