

# Lang Cao

(+1) 217-621-0532 | [lgcao2000@163.com](mailto:lgcao2000@163.com)  
[windszzlang.github.io](https://windszzlang.github.io)

## Education

### University of Illinois at Urbana-Champaign (UIUC)

Urbana, USA

Master of Science in Computer Science (*Research-based Program*)

Sept. 2022 - May 2024 (expected)

- GPA: 3.89/4.0

### Wuhan University of Technology (WHUT)

Wuhan, China

Bachelor of Engineering in Software Engineering

Sept. 2018 - June 2022

- GPA: 4.351/5.0 (WES 3.94/4.0), Rank: 1/79

## Publications

### AutoAM: An End-To-End Neural Model for Automatic and Universal Argument Mining [\[Paper\]](#)

- Lang Cao
- In 19th anniversary of the International Conference on Advanced Data Mining and Applications, ADMA 2023

### CBCP: A Method of Causality Extraction from Unstructured Financial Text [\[Paper\]](#) [\[Github\]](#)

- Lang Cao, Shihua Zhang and Juxing Chen.
- In 2021 5th International Conference on Natural Language Processing and Information Retrieval, NLPPIR 2021

### Clustering of Functionally Related Genes Using Machine Learning Techniques [\[Paper\]](#)

- Yujing Xue and Lang Cao.
- In 2021 5th International Conference on Compute and Data Analysis, ICCDA 2021.

### Intelligent Cross-sensing Sensor Based on Deep Learning [\[Paper\]](#)

- Lingfei Xu, Jiaming Zhang, Lang Cao, and Xinyu Hu.
- In 2021 6th IEEE International Conference on Signal and Image Processing, ICSIP2021

### PILOT: Legal Case Outcome Prediction with Case Law [\[Paper\]](#)

- Lang Cao, Zifeng Wang, Cao Xiao, Jimeng Sun
- In the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023 (Under Review)

### DiagGPT: An LLM-based Chatbot with Automatic Topic Management for Task-Oriented Dialogue [\[Paper\]](#)

- Lang Cao
- In the 38th annual AAAI Conference on Artificial Intelligence, AAAI 2024 (Under Review)

### Enhancing Reasoning Capabilities of Large Language Models: A Graph-Based Verification Approach [\[Paper\]](#)

- Lang Cao
- In the 38th annual AAAI Conference on Artificial Intelligence, AAAI 2024 (Under Review)

## Experiences

### Sunlab, UIUC

Urbana, US

Research Assistant

Jan. 2023 - Now

- Research Focus: Natural Language Processing for Legal and Healthcare.
- Advisor: Jimeng Sun, Danica Xiao

### LegalDAO (Legal AI Startup)

Beijing, China

NLP Tech Leader @ LegalNOW

July 2023 - Now

- Job Content: Completed the construction of an AI legal chatbot which can help and guide user to achieve multiple functions related to contracts.

### Bioinformatics Innovation Lab (Text Group), WHUT

Wuhan, China

Research Assistant

Jan. 2021 - May. 2022

- Research Focus: Information Extraction and Data Mining.
- Advisor: Jing Peng

### iFLYTEK CO. LTD.

Hefei, China

Lang Cao's Curriculum Vitae, 1

- Job Content: maintained and developed an automatic data iteration algorithm for training data of smart car AI system.
- Advisor: Shen'an Li

## Research Projects

### LLM-based Chatbot with Automatic Topic Management for Task-Oriented Dialogue

Beijing, China

Independent Study

June 2023 – Aug. 2023

- Current LLMs have shown proficiency in answering general questions. However, basic question-answering dialogue often falls short in complex diagnostic scenarios, such as legal or medical consultations.
- These scenarios typically necessitate Task-Oriented Dialogue (TOD), wherein an AI chat agent needs to proactively pose questions and guide users towards specific task completion. Previous fine-tuning models have underperformed in TOD, and current LLMs do not inherently possess this capability.
- We introduce DiagGPT (Dialogue in Diagnosis GPT), an innovative method that extends LLMs to TOD scenarios. Our experiments reveal that DiagGPT exhibits outstanding performance in conducting TOD with users, demonstrating its potential for practical applications.

### Legal Case Outcome Prediction with Case Law

Urbana, USA

Advised by Danica Xiao and Jimeng Sun

June 2023 – Aug. 2023

- Machine learning shows promise in predicting the outcome of legal cases, but most research has concentrated on civil law cases rather than case law systems.
- We identified two unique challenges in making legal case outcome predictions with case law. First, it is crucial to identify relevant precedent cases that serve as fundamental evidence for judges during decision-making. Second, it is necessary to consider the evolution of legal principles over time, as early cases may adhere to different legal contexts.
- We proposed a new model named PILOT (Predicting Legal case Outcome) for case outcome prediction. It comprises two modules for relevant case retrieval and temporal pattern handling, respectively.
- To benchmark the performance of existing legal case outcome prediction models, we curated a dataset from a large-scale case law database. We demonstrate the importance of accurately identifying precedent cases and mitigating the temporal shift when making predictions for case law, as our method shows a significant improvement over the prior methods that focus on civil law case outcome predictions.

### Graph-Based Verification to Enhance Reasoning Capabilities of Large Language Models

Urbana, USA

Independent Study

Oct. 2023 - Dec. 2023

- Large Language Models (LLMs) have showcased impressive reasoning capabilities, particularly when guided by specifically designed prompts in complex reasoning tasks such as math word problems.
- These models typically solve tasks using a chain-of-thought approach, which not only bolsters their reasoning abilities but also provides valuable insights into their problem-solving process. However, there is still significant room for enhancing the reasoning abilities of LLMs. Some studies suggest that the integration of an LLM output verifier can boost reasoning accuracy without necessitating additional model training.
- We follow these studies and introduce a novel graph-based method to further augment the reasoning capabilities of LLMs. We posit that multiple solutions to a reasoning task, generated by an LLM, can be represented as a reasoning graph due to the logical connections between intermediate steps from different reasoning paths.
- We propose the Reasoning Graph Verifier (RGV) to analyze and verify the solutions generated by LLMs. By evaluating these graphs, models can yield more accurate and reliable results. Our experimental results show that our graph-based verification method not only significantly enhances the reasoning abilities of LLMs but also outperforms existing verifier methods in terms of improving these models' reasoning performance.

### End-to-End Argument Mining on Universal Unstructured Text

Wuhan, China

Research Intern, Advised by Prof. Jing Peng

Apr. 2022 - Dec. 2022

- Argument mining is to analyze argument structure and extract important argument information from unstructured text. An argument mining system can help people automatically gain causal and logical information behind the text. As argumentative corpus gradually increases, like more people begin to argue and debate on social media, argument mining from them is becoming increasingly critical. However, argument mining is still a big challenge in natural language tasks due to its difficulty, and relative techniques are not mature. For example, research on non-tree argument mining needs to be done more. Most works just focus on extracting tree structure argument information. Moreover, current methods cannot accurately describe and capture argument relations and do not predict their types.
- We propose a novel neural model called AutoAM to solve these problems. We first introduce the argument component attention mechanism in our model. It can capture the relevant information between argument components, so our model can better perform argument mining. Our model is a universal end-to-end framework, which can analyze

argument structure without constraints like tree structure and complete three subtasks of argument mining in one model. The experiment results show that our model outperforms the existing works on several metrics in two public datasets.

## Causality Extraction from Unstructured Financial Text

Wuhan, China

Advised by Prof. Jing Peng

Jan. 2021 - May 2021

- Extracting causality information from unstructured natural language text is a challenging problem in natural language processing. However, there are no mature special causality extraction systems. Most people use basic sequence labeling methods, such as BERT-CRF model, to extract causal elements from unstructured text and the results are usually not well. At the same time, there is many causal event relations in the field of finance. If we can extract enormous financial causality, this information will help us better understand the relationships between financial events and build related event evolutionary graphs in the future.
- We propose a causality extraction method for this question, named CBCP (Center word-based BERT-CRF with Pattern extraction), which can directly extract cause elements and effect elements from unstructured text. Compared to BERT-CRF model, our model incorporates the information of center words as prior conditions and performs better in the performance of entity extraction. Moreover, our method combined with pattern can further improve the effect of extracting causality. Then we evaluate our method and compare it to the basic sequence labeling method. We prove that our method performs better than other basic extraction methods on causality extraction tasks in the finance field.

## Honors & Awards

---

- Silver Medal, top 5% in Kaggle Common Lit Readability Prize (2021.8)
- Top 2% in Alibaba Tianchi NLP Chinese Pre-training Model Generalization Ability Challenge (2021.1)
- National Scholarship, WHUT (2020); Merit Student Model Honor, WHUT (2020); First-class Scholarship, WHUT (2021); Merit Student Honor, WHUT (2021); Outstanding Graduate (2022.6); Outstanding Graduation Thesis (2022.6)
- China National Championship in the FIRST LEGO League (2014); Asia-Pacific Gold Award in the FIRST LEGO League (2016)

## Skills

---

- **Programming:** Python, C/C++, Java, JavaScript
- **Techniques:** PyTorch, TensorFlow, Huggingface Transformers, LangChain, DGL, Scikit-learn, NumPy, Pandas, Django, Flask
- **Others:** LaTeX, Markdown, Git, SQL, Linux