

# Junjie Chen

[✉ junjiechen.chris@outlook.com](mailto:junjiechen.chris@outlook.com) | [🏡 https://chrisjhc.com/](https://chrisjhc.com/) | [👤 junjiechen-chris](https://www.linkedin.com/in/junjie-chen-b6a075240) | [🔗 junjie-chen-b6a075240](https://www.linkedin.com/in/junjie-chen-b6a075240)

## Summary

---

CS PhD with 5 years of experience in NLP. Specializing in LLM-based unsupervised parsing to extract natural language information from large scale data. Broad knowledge of AI/ML methodologies. Worked on LLM-based book translation and LLM Sparse AutoEncoder-based personal information identification.

## Education

---

### The University of Tokyo, Advised by Yusuke Miyao

PH.D. IN COMPUTER SCIENCE

- Supported by Japan Society for the Promotion of Science (JSPS) Fellowship
- PhD research published in 3 top-tier venues

Tokyo, Japan

Oct. 2021 - Mar. 2025

### The University of Tokyo, Advised by Yusuke Miyao

M.S. IN COMPUTER SCIENCE

- Master research published in ACL 2022 (Oral)

Tokyo, Japan

Oct. 2019 - Sep. 2021

### University of Liverpool, Advised by Danushka Bollegala

VISITING SCHOLAR

- Research published in ACL 2024 Findings

Liverpool, UK

Oct. 2023 - Mar. 2024

### Northeastern University

B.S. IN SOFTWARE ENGINEERING

- GPA: 3.46/4.0 (3.78/4.0 during 6 month exchange at Hokkaido University)

Shenyang, China

Sep. 2015 - Jun. 2019

## Selected Publications

---

### Improving Unsupervised Constituency Parsing via Maximizing Semantic Information

ICLR 2025 (Spotlight, top 5.1%)

JUNJIE CHEN, XIANGHENG HE, YUSUKE MIYAO, AND DANUSHKA BOLLEGALA

2025

### Constituents are Frequent Word Sequences among Sentences with Equivalent Predicate-Argument Structures: Unsupervised Constituency Parsing by Span Matching

ACL 2024 Findings

JUNJIE CHEN, XIANGHENG HE, DANUSHKA BOLLEGALA, AND YUSUKE MIYAO

2024

### Language Model Based Unsupervised Dependency Parsing with Conditional Mutual Information and Grammatical Constraints

NAACL 2024

JUNJIE CHEN, XIANGHENG HE, AND YUSUKE MIYAO

2024

### Syntactic-Semantic Dependency Correlation in Semantic Role Labeling: A Shift in Semantic Label Distributions

Journal of Natural Language Processing

JUNJIE CHEN

2022

### Modeling Syntactic-Semantic Dependency Correlations in Semantic Role Labeling Using Mixture Models

ACL 2022 (Oral)

JUNJIE CHEN, XIANGHENG HE, AND YUSUKE MIYAO

2022

### ProsodyFM: Unsupervised Phrasing and Intonation Control for Intelligible Speech Synthesis

AAAI 2025 (Oral)

XIANGHENG HE, JUNJIE CHEN, ZIXING ZHANG, BJÖRN W SCHULLER

2025

### Task Selection and Assignment for Multi-modal Multi-task Dialogue Act Classification with Non-stationary Multi-armed Bandits

ICASSP 2024

XIANGHENG HE, JUNJIE CHEN AND BJORN W. SCHULLER

2024

**A system for worldwide COVID-19 information aggregation**

Workshop on NLP for COVID-19

AKIKO AIZAWA, FREDERIC BERGERON, JUNJIE CHEN, FEI CHENG, KATSUHIKO HAYASHI, KENTARO INUI, HIROYOSHI ITO,

DAISUKE KAWAHARA, MASARU KITSUREGAWA, HIROKAZU KIYOMARU, MASAKI KOBAYASHI, TAKASHI KODAMA, SADAO

KUROHASHI, QIANYING LIU, MASAKI MATSUBARA, YUSUKE MIYAO, ATSUYUKI MORISHIMA, YUGO MURAWAKI, KAZUMASA

OMURA, HAIYUE SONG, EIICHIRO SUMITA, SHINJI SUZUKI, RIBEKA TANAKA, YU TANAKA, MASASHI TOYODA, NOBUHIRO UEDA,

HONAI UEOKA, MASAO UTIYAMA, YING ZHONG

2020

**A pattern-based method for medical entity recognition from Chinese diagnostic imaging text**

Frontiers in Artificial Intelligence

ZIHONG LIANG, JUNJIE CHEN, ZHAOPENG XU, YUYANG CHEN, TIANYONG HAO

2019

**A Bibliometric Analysis of the Research Status of the Technology Enhanced Language Learning**

Emerging Technologies for Education

XIELING CHEN, JUNTAO HAO, JUNJIE CHEN, SONGSHOU HUA, TIANYONG HAO

2018

## Work Experience

---

### Amazon

Tokyo, Japan

AMAZON SCIENCE FELLOW

Aug. 2025 - Feb. 2026 (Est.)

- Early member and core contributor of a multimodal book translation project
- Designed and implemented a human evaluation process.
- Designed and implemented a LLM-centric AI translation pipeline. Work including text translation, image segmentation, and auto-eval.

### Rakuten

Tokyo, Japan

RESEARCH INTERN

Nov. 2024 - Feb. 2025

- Investigated the business application of LLM Sparse AutoEncoder (SAE)
- Improved SAE-based Personal Information Identification (PII) system from 72% to 87%.

### National Institute of Informatics

Tokyo, Japan

RESEARCH ASSISTANT

Jun. 2024 - Feb. 2025

- Developed evaluation pipeline for Japanese LLM of various sizes (180M to 175B)
- Maintained 800-GPUs clusters for LLM training

### The University of Tokyo

Tokyo, Japan

TEACHING ASSISTANT

Oct. 2022 - Feb. 2023

- Mentored 3 student teams on their final NLP projects (one published at JNLP-2023 conference)

## Research Projects

---

### Amazon

WHOLE-BOOK TRANSLATION PROJECT

2024-2025

- Developed an MQM-style guideline tailored to evaluating book translation.
- Developed an LLM-centric translation system that jointly performs visual grounding, speech bubble detection, story comprehension, and translation.
- Improved annotation consistency from 0 (no agreement) to 0.3 (fair agreement) and translation quality from non-publishable (MQM score 80) to near-publishable (MQM score 95)
- Improved the fluency and story accuracy of the translation.
- Improved the performance of speech bubble detection.

### University of Tokyo, PhD Research Project - 2

UNSUPERVISED CONSTITUENCY PARSING VIA DISCOVERING SEMANTIC INFORMATION

2024-2025

- Developed a paraphrasing-based bag-of-substring model to estimate substring-level semantic information
- Developed a calibrated training objective for PCFG-based unsupervised constituency parsing
- The training objective significantly improved parsing accuracy by 8% absolute on average across four languages
- The semantic information estimate is effective in identifying semantic arguments

## University of Tokyo, PhD Research Project - 1

ENHANCING LANGUAGE MODEL-BASED SYNTACTIC DEPENDENCY PARSING WITH GRAMMATICAL PRIOR

2023-2024

- Developed a Metropolis-Hastings based algorithm for estimating word-word mutual information
- Developed a sub-sampling method to incorporate grammatical prior into mutual information calculation
- The grammatical prior achieved aggregated accuracy improvement by over 5%
- The prior enabled over 10% accuracy uplift in predicting semantic-related dependencies

## University of Tokyo, Master Research Project

ENHANCING SYNTAX-AIDED SEMANTIC DEPENDENCY PARSING WITH MIXTURE-OF-EXPERT MODELS

2021-2022

- Developed a mixture-of-experts models that model semantic dependencies by their syntactic patterns
- Applied variational inference techniques to automatically cluster similar syntactic patterns
- Significantly improved parsing accuracy compared to state-of-the-art baselines

## University of Tokyo, Extracurricular Project

EXTRACURRICULAR PROJECTS

2024 - 2025

- Improved Japanese glyph generation with paired positive-negative prompt engineering
- Improved character-level language modeling with a novel distance-to-whitespace loss function

## Grants & Fellowships

---

2023 **DC2 Fellowship**, Japan Society for the Promotion of Science  
2024 **Special Allowance for Outstanding Student**, Japan Society for the Promotion of Science  
2025 **Travel Grant (\$2000)**, Association for the Advancement of Artificial Intelligence  
2022 **IST-RA Fellowship**, The University of Tokyo

## Skills

---

**Research Experience** Unsupervised Natural Language Parsing with LLM, MultiModal book Translation

**Generative Modeling** Language Models, Diffusion, Normalizing Flow, Variational AutoEncoder

**Machine Learning Frameworks** Pytorch, Tensorflow, Lightning, Hydra, vLLM

**Distributed Computing** Slurm, Huggingface Accelerate, Megatron, Lightning

**DevOps** AWS, Docker, Cloudflare, Brazil, Docker

**Programming Skills** Python, Typescript, Bash

**Language** Chinese (Native), Japanese (N1, Business), English (C-1 equivalent, Business Level)