

# Lan Sang

lan.sang@colorado.edu | 2300 Arapahoe Ave, Boulder, Colorado, USA

## EDUCATION

<b>University of Colorado Boulder</b> Ph.D. in Information Systems	Boulder, CO, United States Expected 2026
<b>University of Colorado Boulder</b> M.S. in Computational Linguistics	Boulder, CO, United States May.2020
<b>Nanjing Normal University</b> B.A. in French and BBA in Business Administration	Nanjing, China Jun. 2018

## PROFESSIONAL EXPERIENCE

<b>Nuance Communications, Inc.</b> <i>NLP Research Intern, Natural Language Understanding Team</i>	Boulder, CO, United States Jun. 2019-Aug. 2019
<ul style="list-style-type: none"><li>• Mapped annotated data from a Mandarin corpus in one annotation spec to another annotation spec</li><li>• Ran mapping rules of English data on Mandarin data to see how the rules performed and found out possible reasons that affected the accuracy of the mapping rules when ran on another language</li><li>• Presented to the team about key findings by the end of the internship</li></ul>	
<b>University of Colorado Boulder</b> <i>Research Assistant, Computational Language and Education Research Center</i>	Boulder, CO, United States May. 2019-Jan. 2021
<ul style="list-style-type: none"><li>• Used html, python (flask framework) and MongoDB to build and maintain the website of United Verbs Index (uvi.colorado.edu), which is a system merges links and web pages from five different NLP projects: VerbNet, FrameNet, PropBank, OntoNotes and SynSemClass Lexicon</li><li>• Edited html and css files to improve the functions and user interface of the UVI website as well as updating VerbNet</li></ul>	

## RESEARCH EXPERIENCE & PROJECTS

<b>SIGMORPHON 2020 Shared Task: Typologically Diverse Morphological Inflection</b> <i>Course Project</i>	Boulder, CO, United States Mar.2020-May.2020
<ul style="list-style-type: none"><li>• Wrote grammar rules to generate lemmas, part-of-speech tags and word forms for three languages: Crimean Tatar (crh), Tagalog (tgl) and Livonian (liv)</li><li>• Built Finite-State Transducer (FST) using Foma by combining a lexicon-based model with a guesser to handle unseen lemmas and got accuracies of 96.38%(crh), 68.84%(liv) and 78.01%(tgl) on test datasets.</li></ul>	
<b>Topic Modeling Using Unsupervised Learning Models on User Review Dataset</b> <i>Course Project</i>	Boulder, CO, United States Nov.2019
<ul style="list-style-type: none"><li>• Clustered customer reviews into groups and discovered the latent semantic structures using Python</li><li>• Preprocessed text dataset of Amazon user reviews by tokenization, stemming and extracted features by Term Frequency-Inverse Document Frequency (TF-IDF)</li><li>• Trained unsupervised machine learning models of K-Means Clustering and Latent Dirichlet Allocation (LDA)</li><li>• Visualized and analyzed the model training results</li></ul>	

## SKILLS

- Background: Computational Linguistics | Natural Language Processing | Machine Learning | Deep Learning
- Programming: Python | Java | SQL | JavaScript
- Tools: MySQL | Spark | Tensorflow | Keras | PyTorch | AWS | Tableau | Git | Praat | NLTK | Scikit-learn | NumPy
- Languages: English (Fluent) | Mandarin (Native) | French (Fluent)

## PUBLICATIONS

- Sarah Beemer, Zak Boston, April Bukoski, Daniel Chen, **Lan Sang**, Mans Hulden et al(2020), "Linguist vs. Machine: Rapid Development of Finite-State Morphological Grammars." In *Proceedings of 17<sup>th</sup> SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, Pages 162 – 170, Association for Computational Linguistics