# ontoNLP: A Mechanism for Transforming Textual Data into Structured Data

Jun WANG[1,2],  Tong WEI[1,2,*],  Xuemei TANG[1,2],  Zhaoji WANG[3]

1. The Department of Information Management, Peking University
2. The Research Center of Digital Humanities, Peking University
3. School of International Chinese Language Education, Beijing Normal University

**Abstract** In the realm of digital humanities, the focus is on steering humanities research towards data-driven methodologies. Data-driven humanities research commonly relies on structured datasets. However, for humanities scholars, effectively converting textual data into structured data poses a significant challenge. This challenge is especially pronounced when dealing with large-scale textual data under unified standards. Annotation plays a pivotal role in dealing with this challenge. As a response, this article delves into the intersection of ontology and natural language processing (NLP) technology within the realm of annotation. It proposes an ontoNLP mechanism for automated annotation based on pre-defined ontology, aiming to facilitate the transformation of textual data into structured datasets. And, we developed an online system WuYuDian (吾与点) for automatic tagging ancient Chinese texts which is now openly accessed on Web.

## 1. Introduction

Ontology, a cornerstone in the field of knowledge representation, serves as a structured framework for capturing and organizing concepts and relationships within a specific domain. At its core, ontology provides a shared vocabulary and a set of defined rules that enhance communication and understanding of concepts among humans, and crucially, between humans and machines. Natural Language Processing (NLP) stands at the crossroads of linguistics and artificial intelligence, with the aim of bridging the gap between human language and computer comprehension.

Annotation is the core task of transforming textual data into structured data. OntoNLP is a mechanism designed to convert textual data into structured data. Its fundamental idea is to harness the strengths of both ontology and NLP technology through annotation, facilitating to process large-scale text data and organize knowledge in a unified annotation scheme. Based on ontoNLP mechanism, we developed an online automatic annotation tool —— 吾与点 (Wu-Yu-Dian). WYD is designed to facilitate the training data generation for machine learning, and the data extraction from textual materials for data-driven digital humanities works.

## 2. ontoNLP mechanism

ontoNLP mechanism is guided by the research questions, which is different from other annotation tools (Figure 1). In this mechanism, users first identify the research question, then collect textual data and define the domain ontology based on the research question. The concept layer in the domain ontology guides named entity annotation, while relationships guide relationship

annotation. After the annotation is completed, manual proofreading and modification are required, followed by the export of structured data. In this mechanism, ontology, as a formal and clear representation of knowledge and a conceptual framework for understanding the relationships between entities, is the cornerstone of automatic annotation. It provides a structured vocabulary and a set of rules to guide the extraction of information from unstructured text. NLP, when combined with ontology, helps extract contextual information from text and achieve automatic annotation of large-scale texts.
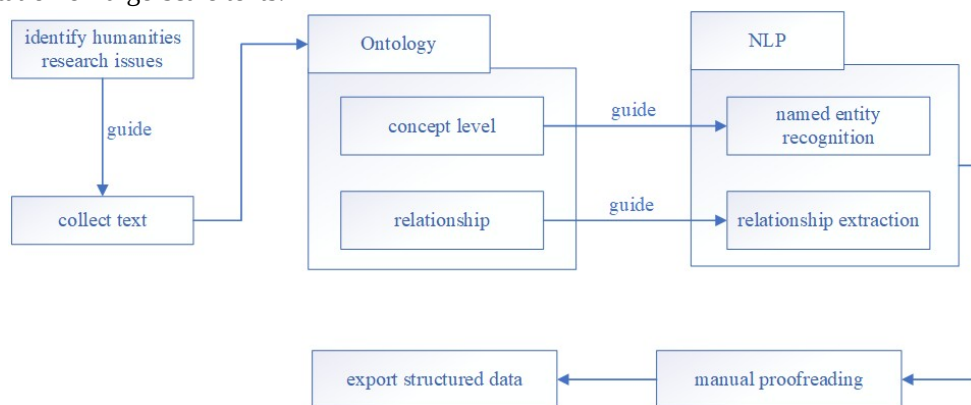


Figure 1. ontoNLP mechanism.

## 3. WuYuDian: An online annotation tool based on ontoNLP

Based on the ontoNLP mechanism, we developed WuYuDian that is an online automatic annotation tool (Figure 2) capable of annotating entities and extracting relationships. It is accessible at: https://www.wyd.pkudh.org. This tool can currently only handle Chinese ancient texts, and we will expand it into a multilingual processing tool in the future. The core functions of this tool include entity annotation (Figure 2) and relationship annotation (Figure 3). Entity annotation comprises concept layer, individual layer, and text layer. Users first define a conceptual model and upload it at the concept layer. Subsequently, they can initiate automatic annotation by clicking the auto-annotation button based on their conceptual model. Upon completing automatic annotation, users can manually proofread and modify in the text area. The automatic annotation function is facilitated through named entity recognition, achieved by transfer learning and incremental training on the BERT model. Currently, the entities that can be automatically annotated encompass five types: time, location, person, official, and book title, with an accuracy rate of approximately 92% in ancient Chinese texts. Furthermore, the platform offers users manual annotation and proofreading functions, with annotated entity results summarized and displayed at the individual level.

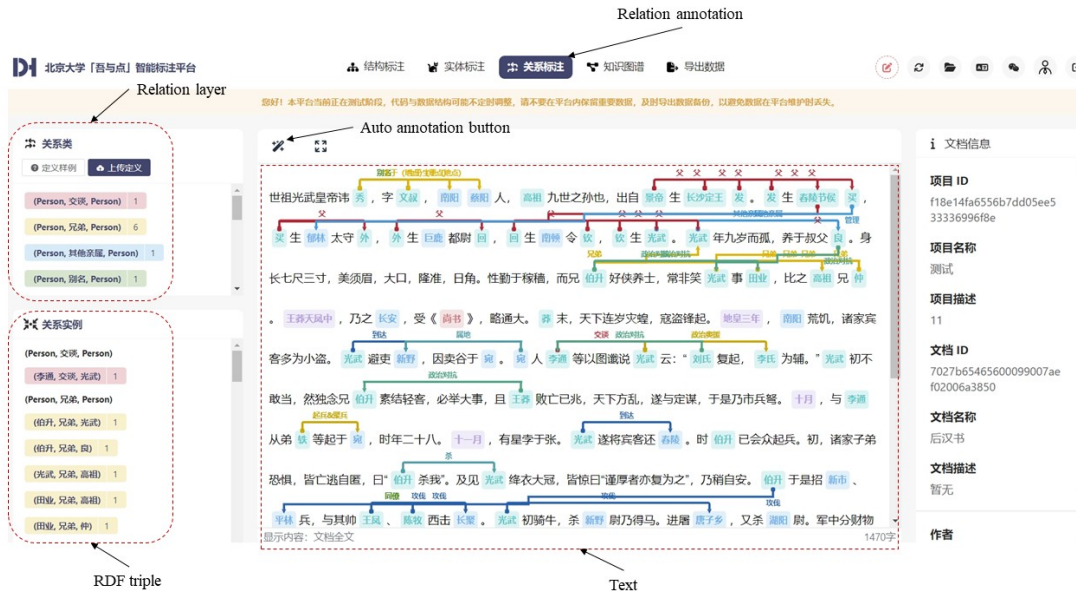Figure 2. WuYuDian: an online tool of automatic annotation.



Figure 3. The relationship annotation

Relationship annotation (Figure 3) also comprises three parts: the relationship layer, RDF triplet, and text layer. Users are required to upload the defined relationships in the relationship layer and then initiate the automatic annotation by clicking the auto-annotation button, which automatically annotates the relationships in the text layer. Relationship annotation is based on the entities identified during the entity annotation process. Once users finish the automatic annotation, they can manually proofread and modify in the text area. The relationship annotation currently supports nearly 30 types of relationships, including aliases, appointments, political alliances, conquests, management, familial relationships, etc. The model achieves an 80% accuracy in extracting relationships from historical texts.

## 4. Conclusion and Future work

This article primarily introduces the ontoNLP mechanism based on ontology and NLP technology for converting textual data into structured data. It then presents an online data annotation tool

developed using the ontoNLP mechanism, focusing on entity annotation and relationship annotation functions. Annotation is a critical step in converting textual data into structured data. The WuYuDian online annotation tool developed in this article aims to assist humanities scholars in constructing datasets and conducting data-driven research. However, a limitation of this approach is its reliance on the NLP model, resulting in the inability to automatically recognize new entity types. To address this, our future work will involve training multi-domain and cross-language Bert models to enhance service provision or explore alternatives with large language models.