# Long Legal Document Entity Extraction with Hierarchical Section-based Representation

Data Science Institute
COLUMBIA UNIVERSITY

Authors: Ning Kang, Xindi Deng, Yihan Chen, Yuhui Wang, Yuxin Lin

Mentor: Vladimir A. Kobzar

Data Science Capstone Project with Fu Foundation School of Engineering and Applied Science

## Background

Owing to the intricate and lengthy nature, the reading comprehension of legal documents raises great challenges for both human and AI models. To help readers grasp the main idea in the Merger and Acquisition (M&A) Agreements, the MAUD paper proposes 22 types of questions and aims to find key clauses (i.e., answers) for each question. Each question relates to a deal point type and the answer is a span in the contract. However, the original solution in MAUD did not achieve a satisfying result.

**Main idea:** Use section hierarchy information to help model understand long legal document.

## Introduction

**Dataset:** This project uses hierarchical representation of text from MarkupMnA and deal point information from MAUD. These two datasets share the same 151 M&A agreements.
- **MarkupMnA:** CSV formatted contracts containing annotated sections titles and section numbers, derived from HTML.
- **MAUD:** JSON formatted contracts with 22 questions and relevant answers per agreement. Each question corresponds to a deal point type.

**Task:** Extractive question-answering task on M&A agreements as example shown in Figure 1.

**Contract:** ......a material breach of this Agreement; (e) by Parent, if, prior to obtaining the Company Stockholder Approval, (i) the Company Board of Directors shall have effected a Change of Recommendation (whether or not in compliance with this Agreement) or (ii) the Company has materially breachedSection5.3; (f) by either the Company ....... Meeting duly convened therefor or at any adjournment or postponement thereof; or (h) by the Company in order to effect a Change of Recommendation and substantially concurrently enter into a definitive agreement providing for a Superior Proposal

**Question:** Highlight the parts of the text related to deal point "Breach of No Shop" that should be reviewed by a lawyer
**Answer:** (e) by Parent, if, prior to obtaining the Company Stockholder Approval, (i) the Company Board of Directors shall have effected a Change of Recommendation (whether or not in compliance with this Agreement) or (ii) the Company has materially breachedSection5.3;

**Question:** Highlight the parts of the text related to deal point "Lim. on FTR exercise" that should be reviewed by a lawyer
**Answer:** (h) by the Company in order to effect a Change of Recommendation and substantially concurrently enter into a definitive agreement providing for a Superior Proposal

Figure 1. Example of Task

## Data Overview

MarkupMnA dataset is in CSV format, containing hierarchy representations that relates text spans to different levels of hierarchy. With the visualization tool, we found that multiple CSV files had inaccurate hierarchy information.

MAUD dataset is a SQuADv2 formatted JSON file, in which a list of answers to 22 questions are recorded for each agreement.

**Data Source:** Contract in HTML → MarkupMnA → Author-annotated hierarchical spans → CSV Format → Hierarchical span validated with visual highlights / Output section → Section Hierarchy

**Data Source:** Contract in Plain Text → MAUD → Expert-extracted answer spans → JSON Format → Convert JSON to SQuAD examples → Question-Answers Pairs
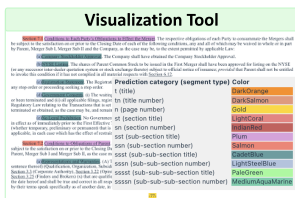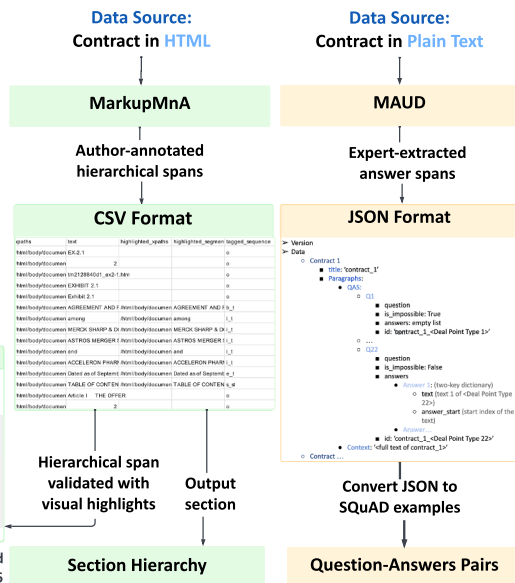
Visualization Tool

Figure 2. Highlighted Hierarchical Text based on MarkupMnA Annotation for Contract 126

## Methods & Models

**Main idea:** To improve the performance of extractive question-answering task on MAUD, we combined section hierarchy information given by MarkupMnA with original solutions proposed by MAUD.

MAUD → Questions / Contract / Answers → Model Input → QA Model → Prediction Spans
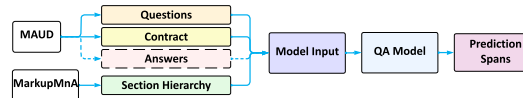MarkupMnA → Section Hierarchy

Figure 3. Workflow

**Baseline Model:** pretrained RoBERTa given by MAUD paper

### Method 1: Section-based Sliding Windows

**Motivation:** For original sliding windows shown in Figure 4, MAUD used a fixed stride length to segment, which causes some subcontracts across different sections, leading to a loss of the contextually coherence and a challenge for model's understanding.

**Contribution:** Incorporating section hierarchy information, we propose a new method to create sliding windows. The paragraph length and stride length are customized with the section split points as Figure 5. If the current sliding window covers a section split point, the split point will be considered as the end of current sliding window and the start of the next sliding window.
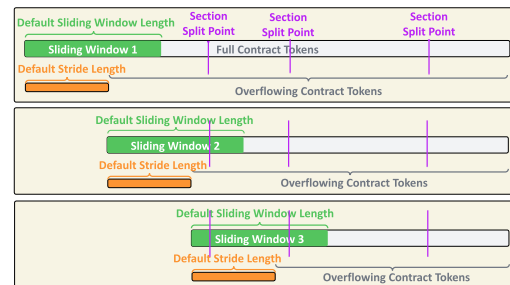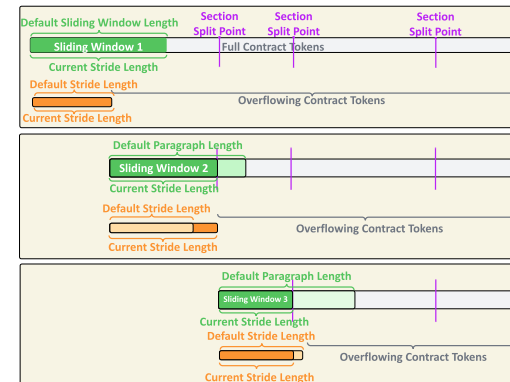
Figure 4. Original Sliding Window

Figure 5. Section-based Sliding Window

### Method 2: Section-based Training Set

**Motivation:** The original model used the data processing functions from Hugging Face. Corresponding to each question-contract pair, there could be many answer spans. However, the functions only extracted the first answer for each question-contract pair. Consequently, the remaining answers are left unused and unexplored.

**Contribution:** We segmented each contract into multiple sections using the section hierarchy mentioned above. In this case, the functions are able to extract one answer from each question-section pair, and more answer spans can be used in total. Therefore, the model can better harness the potential of the dataset.

JSON formatted data: Question: Q / Contract: C / Answers: A1,A2,...,An

Original Examples: Question: Q / Contract: C / Answer: A1

Segmented Examples: Question: Q / Section: S1 / Answers: A1 | Question: Q / Section: S2 / Answers: A2 | ... | Question: Q / Section: Sk / Answers: An

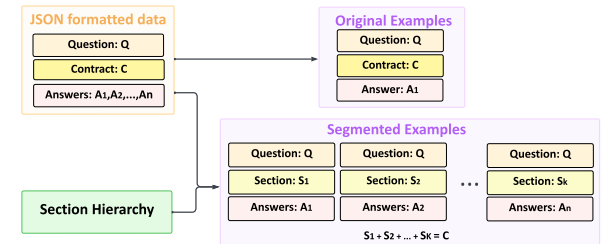$S_1 + S_2 + ... + S_K = C$

Section Hierarchy

Figure 6. Comparison of Original Examples and Segmented Examples

## Results & Conclusion

**Result:** Our Section-based Sliding Window in method 1 outperforms the original method on F1-score for all cases. The fine-tuned model trained on method 2 achieves a higher F1 score for question-contract pairs without cor answers.

**Conclusion:** In spite of the concerns existed in the pretrained model, bringing in the hierarchical information could still make a difference, improving the performance of the extraction task.

**Future Work:** A new embedding of section title could be added to model to take use of global information; The data processing functions could be modified to ensure the use of all the answers.

| Method | F1 score | | | AUPR |
| --- | --- | --- | --- | --- |
| | Total | Has Ans | No Ans | |
| Baseline | 52.59 | 62.76 | 7.69 | 0.2186 |
| Section-based Sliding Window | 54.08 | **63.89** | 10.77 | **0.2205** |
| Section-based Training Set | **55.67** | 51.55 | **73.85** | 0.2136 |

Table 1. F1-score and AUPR

Figure 7. Precision-Recall Curve

Baseline (AUPR=0.2186)
Section-based Sliding Window (AUPR=0.2205)
Section-based Training Set (AUPR=0.2136)
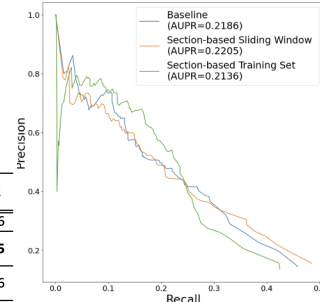
## Acknowledgement

We would like to express our gratitude towards Prof. Kobzar for his support, NYU team for visualization tool template, and Steven H. Wang for the pretrained model.

## References

1. Rao, Sukrit, et al. "MARKUPMNA: Markup-Based Segmentation of M&A Agreements."
2. Wang, Steven H., et al. "MAUD: An Expert-Annotated Legal NLP Dataset for Merger Agreement Understanding." arXiv preprint arXiv:2301.00876 (2023).