



# Decoding Legal Documents: Natural Language Processing and Text Mining for Information Extraction



Irpan Abdurahman, Xinyang Song, Jincal Zheng,  
Lucas Tyler, Tianjun Chen, Aaron Wang

## Introduction

Can you identify the  
“Limit of Liability” value  
from this document?

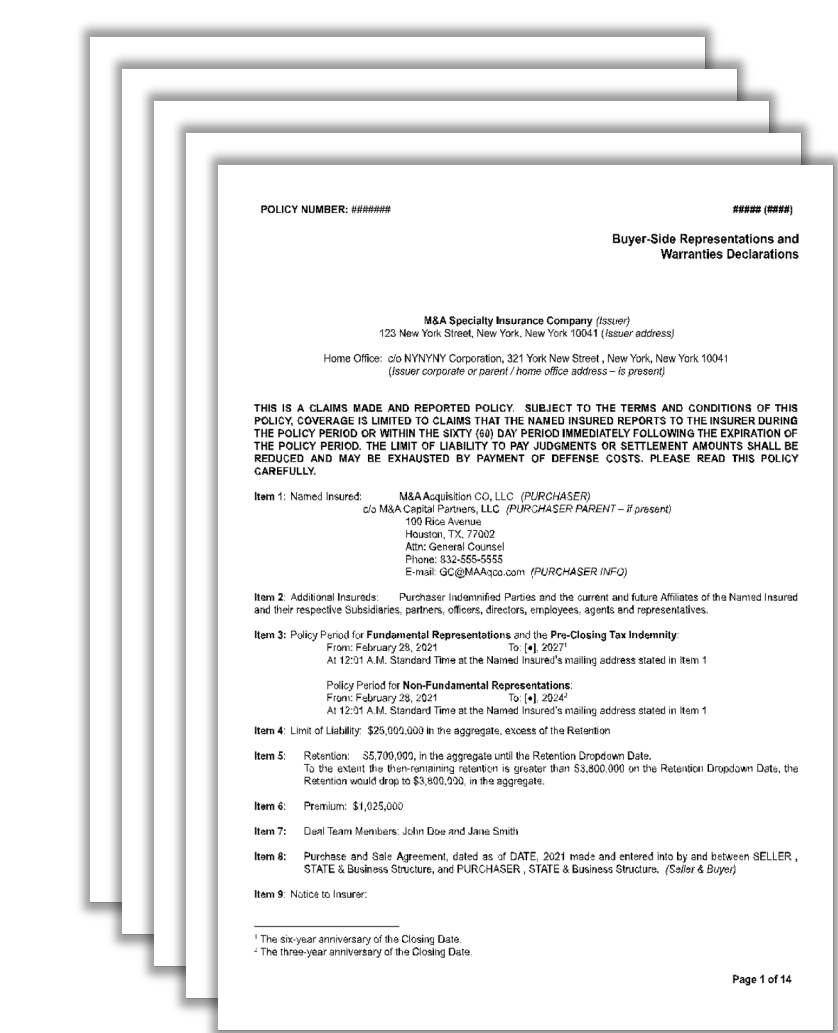
How long did it take?

What if you need to do  
it for 100 documents  
like this?

### Legal Insurance Policy

**Our Solution:** Use Named Entity Recognition and Text Parsing to extract key entities from legal documents.

## Objectives



Doc ID	Premium	...	Policy Ends
001	\$125,000	...	02/28/2024
002	\$800,000	...	06/15/2018
003	\$30,000	...	11/30/2024
004	\$350,000	...	01/05/2025
005	\$700,000	...	12/31/2024

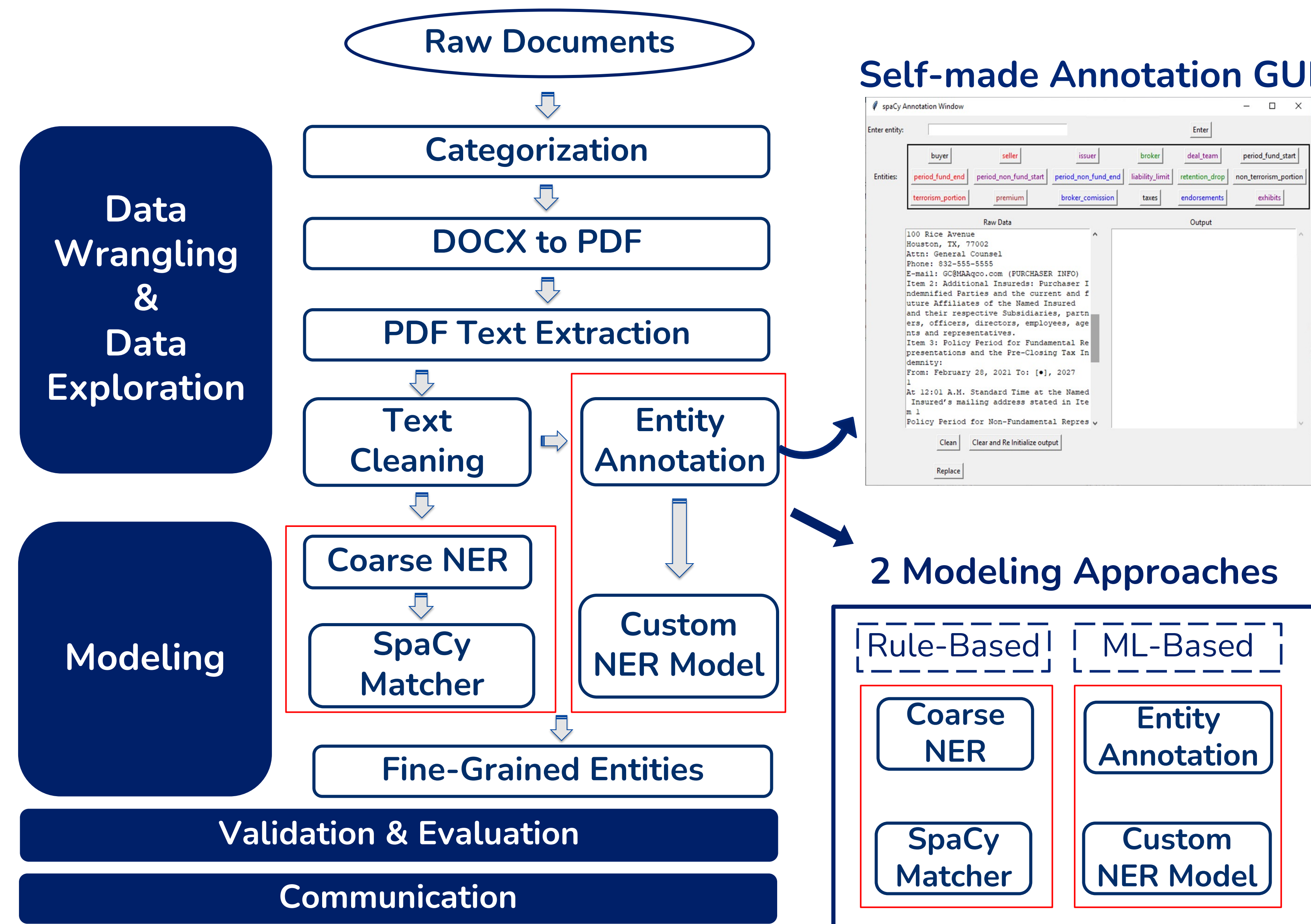
### Structured Database

### Unstructured Legal Documents

Coarse Entity Type	Fine-Grained Entity Type
<b>ID</b>	doc_id
<b>Organization</b>	policy_purchaser issuer broker deal_team
<b>Date</b>	period_fund_start period_fund_end period_non_fund_start period_non_fund_end
<b>Money</b>	premium liability_limit retention retention_drop

### 13 Key Entities to Extract from the Legal Documents

## Pipeline



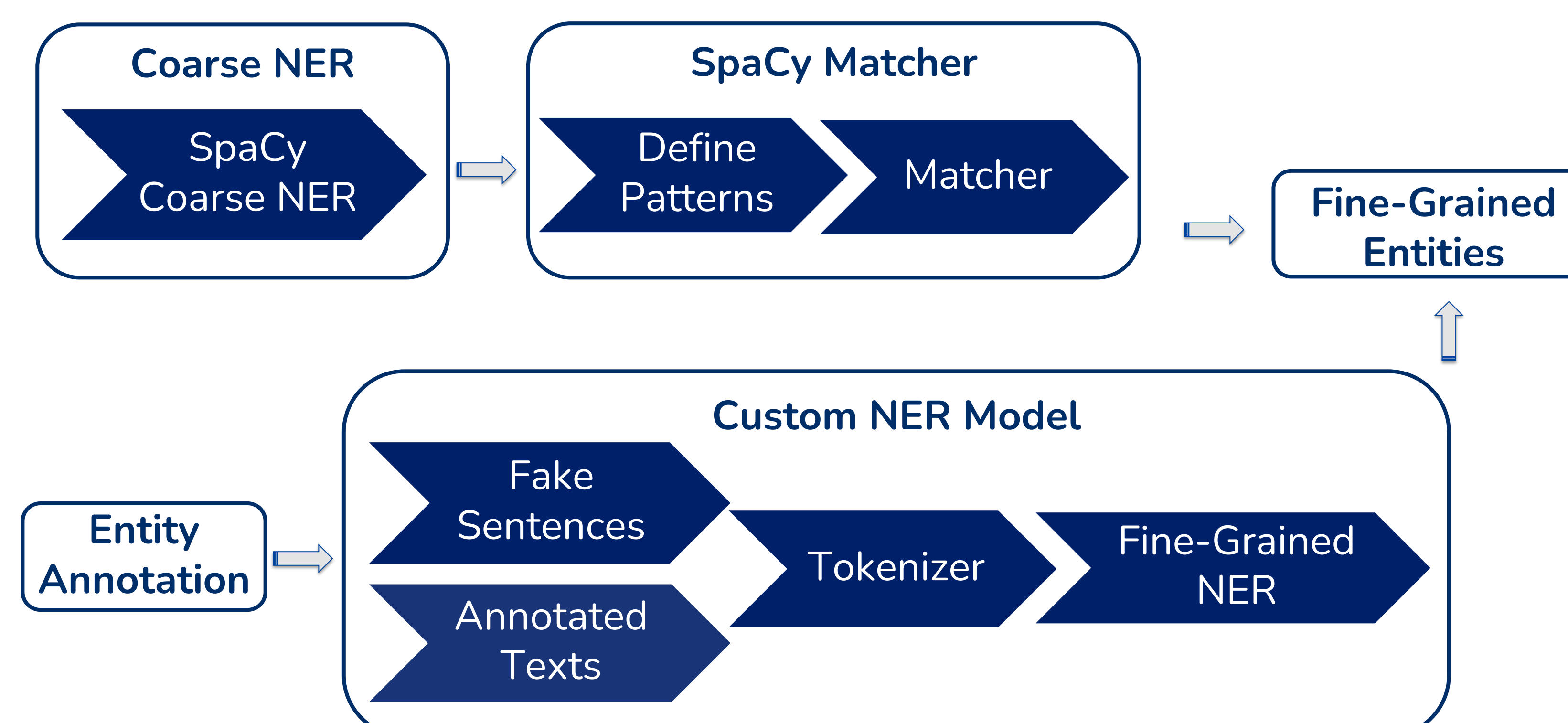
## Methodology

Item 4: Limit of Liability: \$25,000,000 liability\_limit in the aggregate, excess of the Retention (Limit of Liability)

Item 5: Retention: \$5,700,000 retention, in the aggregate until the Retention Dropdown Date.

Item 6: Premium: \$1,025,000 premium

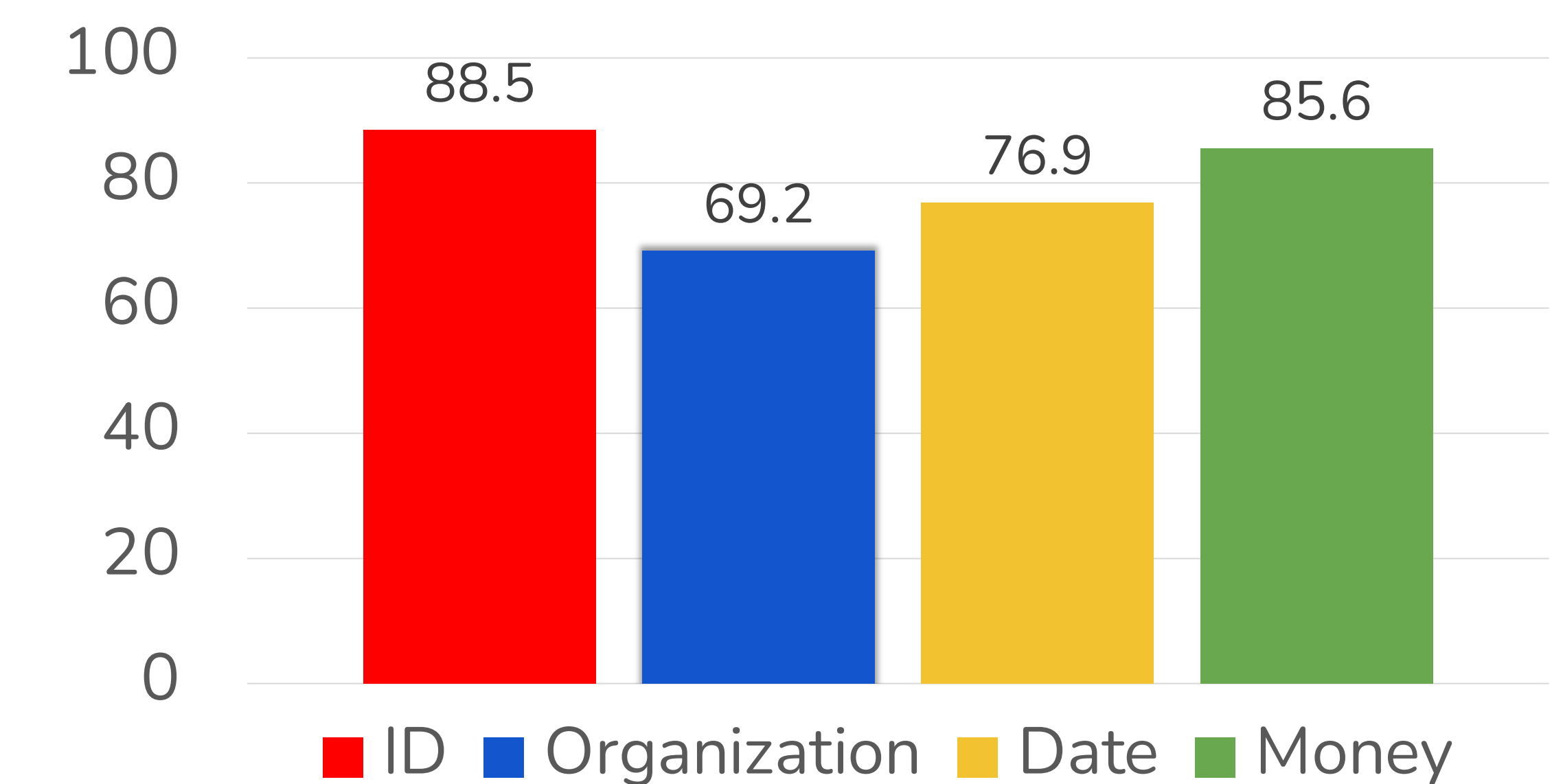
### Named Entity Recognition (NER) Example



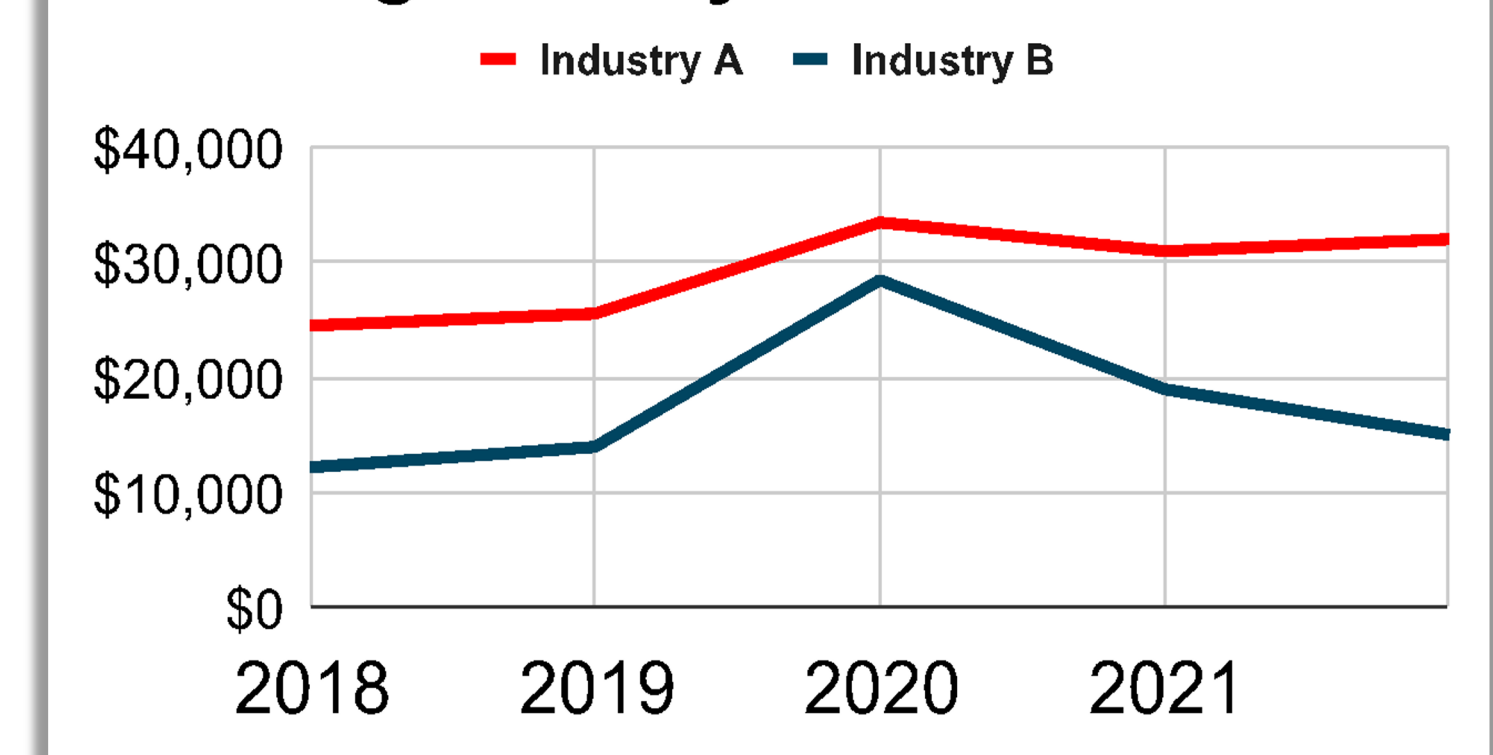
### Rule-based (Top) and ML-based (Bottom) Modeling Methodologies

## Results

### Entity Recognition Accuracy (%)



### Average Policy Premium Value



### Example Use-Case with Extracted Entities

## Conclusion

### Rule-Based

- High accuracy
- Requires maintenance

### ML-Based

- More use-cases
- Need more training data

## Future Work

- Text Extraction**
  - Automate text extraction module
- Rule-Based Model**
  - Use Regex (Time efficiency)
  - More rules for new entity patterns.
- ML-Based Model**
  - More training data (at least 500+ documents).
  - Use the rule-based model to label training data.

## Acknowledgements

Thank you to our sponsors from Norton Rose Fulbright: Zack Barnes, Alper Tuzun, Alexis Collins, and faculty mentor Dr. Arko Barman.