

## Legal Clauses Prediction Model

Prof. Sadia Patka<sup>1</sup>, Mirza Asim Baig<sup>2</sup>, Moin Khan Pathan<sup>3</sup>, Princy Gundetiwar<sup>4</sup>, Adil Ahmed<sup>5</sup>,  
Aaliya Sheikh<sup>6</sup>, Shibani Siddiqui<sup>7</sup>

<sup>1</sup>Assistant Professor, Computer Science and Engineering, A.C.E.T. Nagpur, Maharashtra, India

<sup>2</sup>B.E. Student, Computer Science and Engineering, A.C.E.T. Nagpur, Maharashtra, India

\*\*\*

**Abstract** – There is lack of accessible labelled data that has hindered the application of deep learning in various mastered domains. However, a new solution has been introduced to solve this obstacle within the legal domain of this industry. The creators of this solution are The Atticus Project and they have developed the Contract Understanding Atticus Dataset (CUAD), which is a new dataset for legal contract review. CUAD consists of over 13,000 annotations and was created with the help of dozens of legal experts. The aim of this dataset is to highlight the important portions of a contract that require human verification. Whilst initial results of Transformer models have shown promising performance, there is still a significant room for improvement. This dataset is one of the few large, specialized NLP benchmarks annotated by experts, making it a challenging research benchmark for the wider Natural Language Processing community. Thus, CUAD dataset has the potential to address the bottleneck in deep learning within the legal domain.

**Key Words:** Labelled Datasets, CUAD, Pre-processing, Legal Contract.

### 1. INTRODUCTION

The Contract Understanding Atticus Dataset (CUAD) v1 is a specified dataset consisting of more than 13,000 labels among 510 commercial legal contracts. This dataset was manually labeled by experienced lawyers under the supervision of The Atticus Project, Inc. to identify 41 types of legal clauses that are important for contract review in corporate transactions such as mergers and acquisitions. The CUAD dataset is maintained and curated by The Atticus Project, Inc. to support research and development in Natural Language Processing for legal contract supervision.

Contract review is an essential but a time-consuming work that can cost businesses a good amount of money. For example, many law firms have to spend approximate 50% of their time reviewing contracts, and billing rates for lawyers at large firms can cost between \$500-\$900 per hour in the US. This can result in significant transaction costs for companies, as they have to pay for lawyers to verify that there are no problematic obligations or requirements in the contracts. Furthermore, contract review can be a tedious and monotonous task, which may explain why many individuals and small companies often sign contracts

without reading them. This can result in predatory behavior that harms consumers.

By openly releasing better-quality data (i.e. cleaned data) and models for contract review, improvisation of this task could become possible, and access to legal support for smaller businesses and individuals can exponentially increase. Accessibility of this legal support would be affordable for everyone, rather than just being available to wealthy companies.

### 2. Literature Review

This section includes the related work done in past CUAD dataset.

"CUAD: An Expert-Annotated Natural Language Processing Dataset [1] for Legal Contract Review" presents the Contract Understanding Atticus Dataset (CUAD), which contains over 13k manually annotated labels in 500+ commercial legal contracts. The task is to highlight important legal clauses in the contracts. The study finds that Transformer models have nascent performance and that their performance is influenced by model design and dataset size, leaving room for improvement.

"The Law of Large Documents: Understanding the Format of Legal Contracts Using Visual Cues" proposes a method of segmenting legal contracts based on their structural metadata to improve understanding [2]. This study uses Adam, Attention Dropout, and BERT models and outperforms existing methods on four long-document understanding tasks, as measured on CUAD. This study argues that the segmentation of strategies that are used on longer documents typically miss structural details, leading to poor results on downstream tasks.

"Law Informs Code: A Legal Informatics [3] Approach to Aligning Artificial Intelligence with Humans" describes how data generated by legal processes can facilitate the specification of inherently vague human goals, increasing human-AI alignment and the local usage of AI system. The study proposes a research agenda called "Law Informs Code," which aims to capture complex computational legal processes and embed them in AI, similar to how legal

contracts cannot foresee all potential contingencies, and legislators might not predict all the circumstances under which their proposed bills will be applied.

"Numerical Reasoning over Legal Contracts via Relational Database" [4] study and identify the problem of refining natural language text into a relational database with numerical data structure and querying this database to obtain desired answers. This study however finds that existing approaches are limited in their ability to incorporate domain-specific knowledge and express these mathematical formulas over data structures.

"Con Reader: Exploring Implicit Relations in Contracts for Contract Clause Extraction" proposes a method [5] for automatic Contract Clause Extraction (CCE) by modeling implicit relations in legal contracts. This study analyzes the difficult issues of legal contracts and distills three main implicit relations commonly found in these contracts. Existing CCE methods typically treat contracts as plain text, creating a barrier to understanding complex contracts. This paper introduces a new framework for automatic Contract Clause Extraction (CCE), which takes into account of three important relations: Long-range Context Relation, Term-Definition Relation, and Similar Clause Relation. By exploiting these relations, Con Reader improves contract understanding and CCE performance. The results of the experiments show that Con Reader achieves state-of-the-art performance on two CCE tasks, both in conventional and zero-shot settings. The approach is also found to make predictions more interpretable.

### 3. PROPOSED WORK

This section includes problem statement identified and proposed system.

#### 3.1 PROBLEM STATEMENT

- This legal profession struggles with reviewing very large amounts of legal documents quickly and accurately, which can impede agile and data-driven decision making.
- Reliable and working Artificial Intelligence systems that can understand legal text are urgently needed, but their effectiveness is dependent on the availability, quality, and transparency of annotated legal documents created by legal experts.

Let's understand our main problem with the help of an Example: -

Suppose there is a law student who has to review a bunch of law contract agreements to scan them and extract all the main information such as the document name, name of parties that are involved, the agreement date, the expiration

date etc ... It will take him/her a lot of time to read them and extract all the required information.

#### 3.2 PROPOSED SYSTEM

Progression of the dataset in every next layer also passing through the LSTM model. our project works in a very systematic flow that involves very different layer of which is main layer of our chart ending with the final prediction part which is user based.

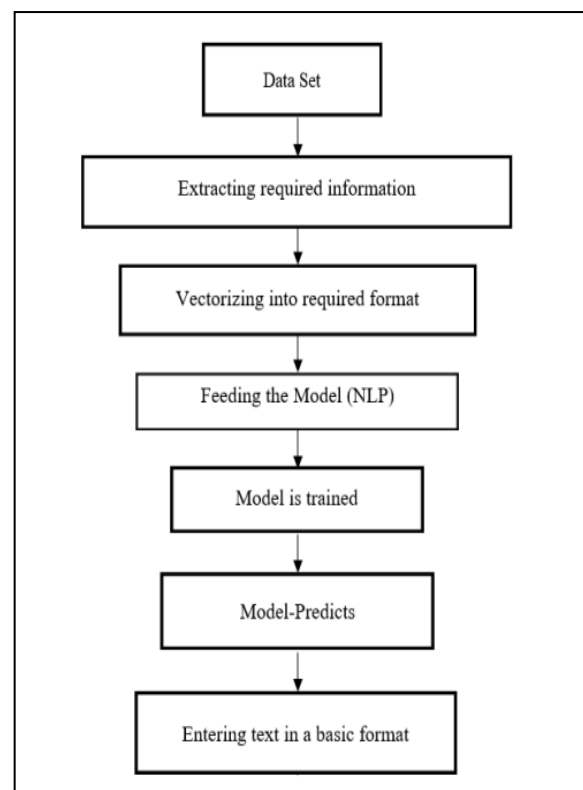


Fig-1 Flow Chart

#### Module 1 : Data Set

In this module, developers collect data about different fruits and then combine the data to create a new data set. Loading

Dataset in Jupyter Firstly, the dataset is downloaded from its website and then is opened in jupyter.

#### Module 2 : Extracting Required Information

Then it is loaded and then some functions will be written in python language to process some parts of the dataset.

### Module 3 : Vectorizing Into Required Format

The data that is processed is then converted into csv file which is then vectorized through the model which contains all the required information

### Module 4 : Feeding The Model

The dataset now is then loaded into the model in the required format in the IOB Tagging method and then is fed into the model.

### Module 5 : Model Is Trained

We have to then load pre trained ELMo model from LSTM keras tensorflow library will which will be used to train the data that we fed into the model ELMo works in such that way that it makes easier to process the language for a computer.

### Module 6 : Model Predicts

Making Predictions Now that the model will be trained, we will use different text format to test the prediction s of the model and then lastly deploy it on a web server using flask.

## 4. RESULT AND DISCUSSION

This is the final IOB tagging function and it's Data Frame created for each clause.

```
Parties_tag()
def Effective_Data_Tag(h):
    answer_1 = data['data'][h]['paragraphs'][0]['cas'][3]['answers']
    if answer_1 == []:
        offset_of = 0
    else:
        res = len([ele for ele in answer_1 if isinstance(ele, dict)])
        j = 0
        while j < len(res)-1:
            it = len(SemanticTokenizer().tokenize(answer_1[j]['text']))
            answer = data['data'][h]['paragraphs'][0]['cas'][3]['answers'][3]['answer_start']
            offset_of_1 = 0
            offset_of_1 = offset_of['start'].tolist()
            if answer in offset_of_1:
                n = offset_of_1.index(answer)
                for i in offset_of['start']:
                    if i in answer:
                        for k in range(0,1):
                            offset_of.at[0+k, 'Tag'] = 'B-Effective Date'
                        else:
                            offset_of.at[0+k, 'Tag'] = 'I-Effective Date'
            j+=1
        j+=1
    Effective_Data_Tag(x)
    offset_of.insert(0, 'Sentences', f'000({x+1})')
    offset_of.drop(['Position', 'start', 'end'], inplace=True, axis=1)
    #offset_of.to_csv(r'C:\Users\VP\Desktop\STHW\DATA\QFT\IOB_of\document.csv', index = False)
    final_of = final_of.append(offset_of)
    final_of.to_csv(r'C:\Users\VP\Desktop\Major Project\documentname.csv', index = False)
    z+=1
```

Fig-2 IOB Tagger Function

The layering of the modelling is shown below which is loaded through keras function . the summary is displayed with 6 layering of models . Input is taken from the user and is passed onto the next layer with the output as it's input and so on . Total Parameters are mentioned as well.

```
In [12]: model_new = tf.keras.models.load_model('model_docparty.h5')
In [13]: model_new.summary()
Model: "model"
Layer (type) Output Shape Param #
-----
input_1 (InputLayer) [(None, 200)] 0
embedding (Embedding) (None, 200, 50) 195200
spatial_dropout1d (SpatialD (None, 200, 50) 0
ropout1d)
bidirectional (Bidirectiona (None, 200, 200) 120800
l)
bidirectional_1 (Bidirectio (None, 200, 200) 240800
ral)
time_distributac (TimeDistr (None, 200, 5) 1005
ibuteac)
-----
Total params: 557,805
Trainable params: 557,805
Non-trainable params: 0
```

Fig-3 Model Summary

This is the training and validation graph which shows exact amount of epochs needed for the model to not over train itself. It shows the accuracy and loss we got after the training is completed.

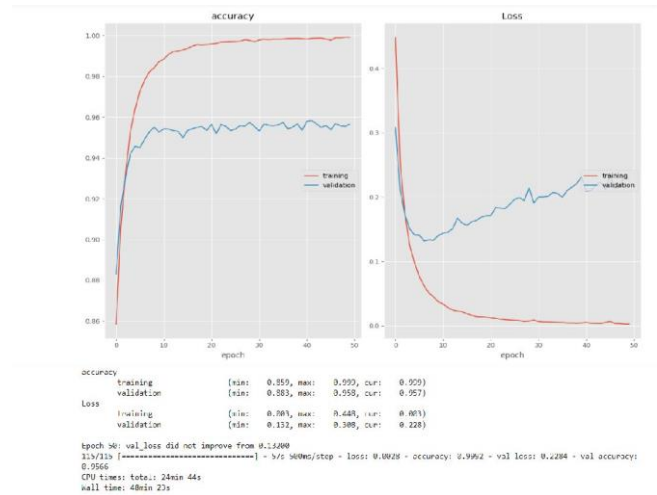


Fig-4 Training and Validation Graph

Finally is the prediction part where the trained model is directly loaded using keras function and then the text is fed into it for the prediction

```
In [28]: i = np.random.randint(0, y_test.shape[0]) #659
p = model.predict(np.array([x_test[i]]))
p = np.argmax(p, axis=-1)
y_true = y_test[i]
print("{}: {} \t {}".format("Word", "True", "Pred"))
print("-" * 30)
for w, true, pred in zip(x_test[i], y_true, p[0]):
    print("{}: {} \t {}".format(words[w-1], tags[true], tags[pred]))
```

Word	True	Pred
S	0	0
SUPPLY	B-Document Name	B-Document Name
AGREEMENT	I-Document Name	I-Document Name
between	0	0
between	0	0
PROFOUND	0	0
MEDICAL	0	0
INC	0	0
.	0	0
and	0	0
PHILIPS	B-Partv	B-Partv

Fig-5 Final Prediction

- [9] <https://www.youtube.com/watch?v=hFUSdgryXyU>
- [10] <https://analyticsindiamag.com/explained-cuad-the-dataset-for-legal-nlp/>
- [11] Suchin Gururangan, Ana Marasovic, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don't stop pretraining: Adapt language models to domains and tasks. ArXiv, abs/2004.10964, 2020.

## 5. CONCLUSION

By far CUAD Dataset is one of the most difficult datasets that can be handled compared to those in the past considering the amount of information that is to be extracted and vectorized into another specific form as required by the LSTM model. The training of the model actually gave a brief idea about how important it is to maintain the accuracy if the model even when changing some of its hyper-parameters. Tuning the model also has a deep effect on the amount of data that is pre-processed into the model. Lastly, the ability to differentiate between the actual pre trained datasets and also the ones which have been improvised.

## REFERENCES

- [1] Deng Cai, Yan Wang, Huayang Li, Wai Lam, and Lemao Liu. 2021. Neural machine translation with monolingual translation memory. In ACL/IJCNLP 2021, pages 7307–7318.
- [2] Ilias Chalkidis and Ion Androutsopoulos. 2017. A deep learning approach to contract element extraction. In JURIX, pages 155–164.
- [3] Afra Nawar, Mohammed Rakib, Salma Abdul Haq. LATERAISSE (LREC) 2022.
- [4] DBAI 2021
- [5] 17 Oct 2022 · Weiwen Xu, Yang Deng, Wenqiang Lei, Wenlong Zhao, Tat-Seng Chua, Wai Lam Fountalis, May.
- [6] Dan Hendrycks, Collin Burns, AnyaChen, SpencrBall CUAD Encodings 10 Mar 2021
- [7] Neel Guha, Daniel E.Ho, Julian Nyarko, Christ opher Rg. Legal Bench. 13 Sep 2022
- [8] Ilias Chalkidis, Ion Androutsopoulos, and Achilleas Michos. 2017. Extracting contract elements. In ICAIL 2017, pages 19–28.