

# Lightweight Contextual Logical Structure Recovery

Po-Wei Huang, Abhinav Ramesh Kashyap, Yanxia Qin,  
Yajing Yang, Min-Yen Kan

*National University of Singapore*

**SDP@COLING 2022**



**NUS**  
National University  
of Singapore

**Computing**



Web Information Retrieval  
Natural Language Processing Group

# Lightweight Contextual Logical Structure Recovery

Po-Wei Huang, Abhinav Ramesh Kashyap, Yanxia Qin, Yajing Yang, Min-Yen Kan\*

National University of Singapore

{huangpowei, abhinav, qinyx, yang0317, kanmy}@comp.nus.edu.sg

## Abstract

Logical structure recovery in scientific documents associates text with a semantic section of the article. Although previous work has disregarded the surrounding context of a line, we model important information by employing line-level attention on top of a transformer-based scientific document processing pipeline. With the addition of loss function engineering and

## Problem Statement

Classify each individual line into 23 predefined classes that indicate the hierarchy of the document structure.

systems (such as Optical Character Recognition (OCR)) to obtain such less cumbersome and similar performance without relying on

by creating a parsimonious model that operates on purely 2D features. Such features

Can we obtain (near-)SOTA  
performance on logical structure  
recovery without relying on feature-  
rich information, but on context only?

# Dataset

## *Data Source*

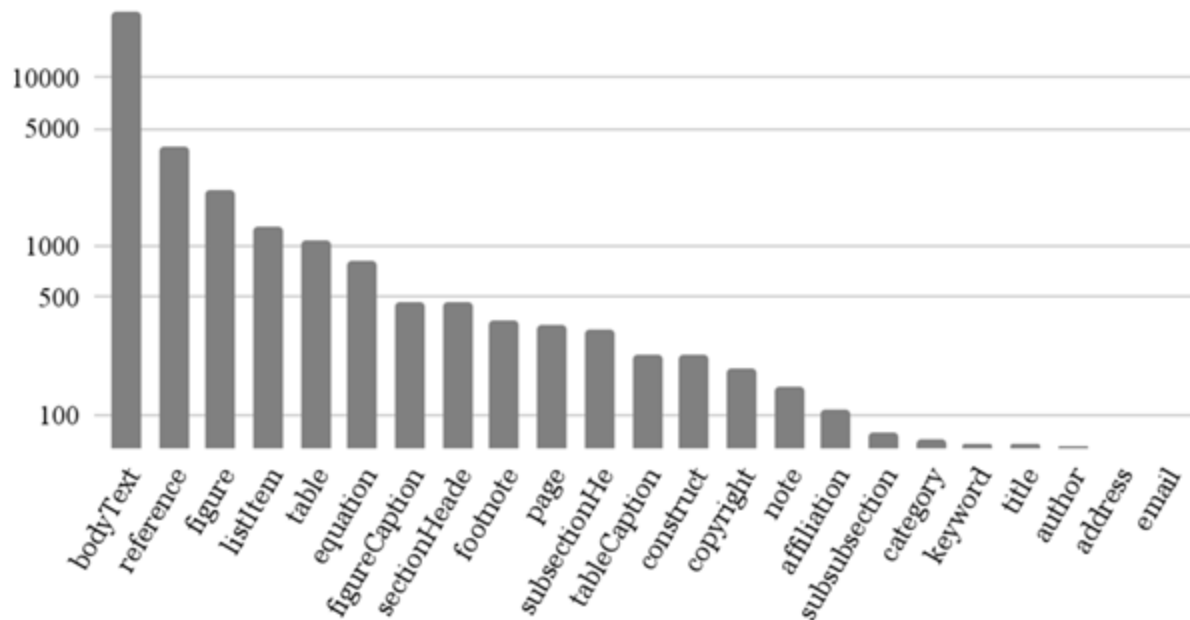
- Original SectLabel Dataset (*Luong et al., IJDLIS 2010*):
  - 20 ACL 2009 Papers
  - 20 CHI 2008 Papers
- Extended Testing Dataset:
  - 20 ACL 2020 Papers
- Unlabelled Dataset:
  - 570 ACL 2021 Long Papers
  - 1895 NeurIPS 2021 Papers

8:1:1 Document Split on SectLabel Dataset for Training, Validation and Testing

# Dataset

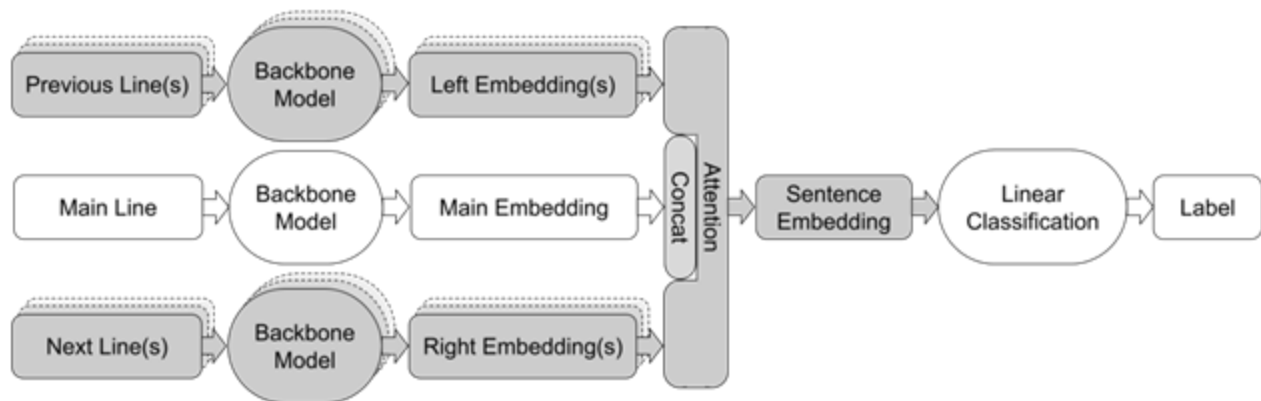
## *Category Distribution*

Occurrence of Each Category



# Contextual Model Construction

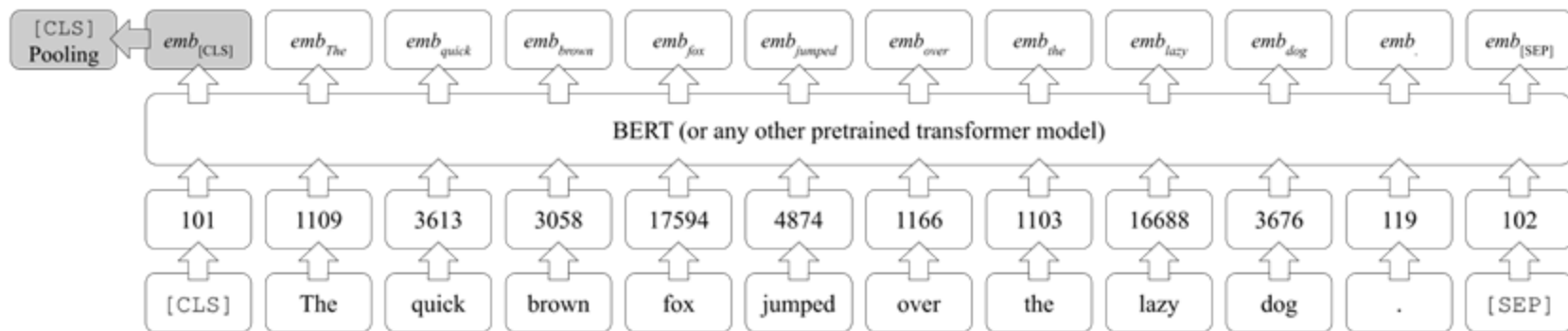
## Sliding Window Attention



	Baseline	Sliding Window 5
Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. arXiv preprint arXiv:1802.05365, 2018.	<i>author</i> reference <i>bodyText</i> reference	reference reference reference reference

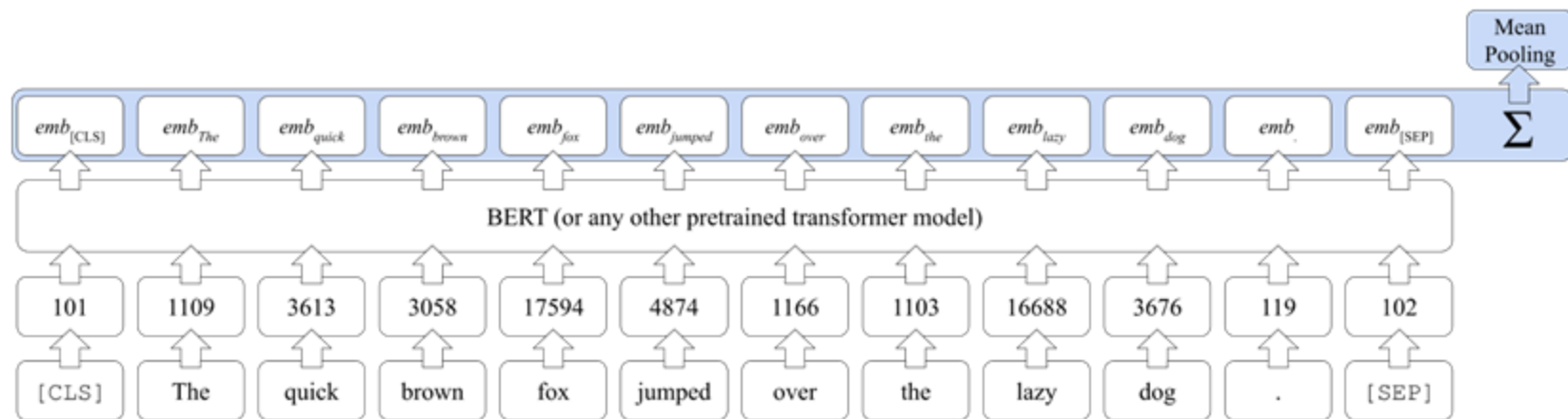
# Pooling Methods for Sentence Embeddings

[CLS] Token



# Pooling Methods for Sentence Embeddings

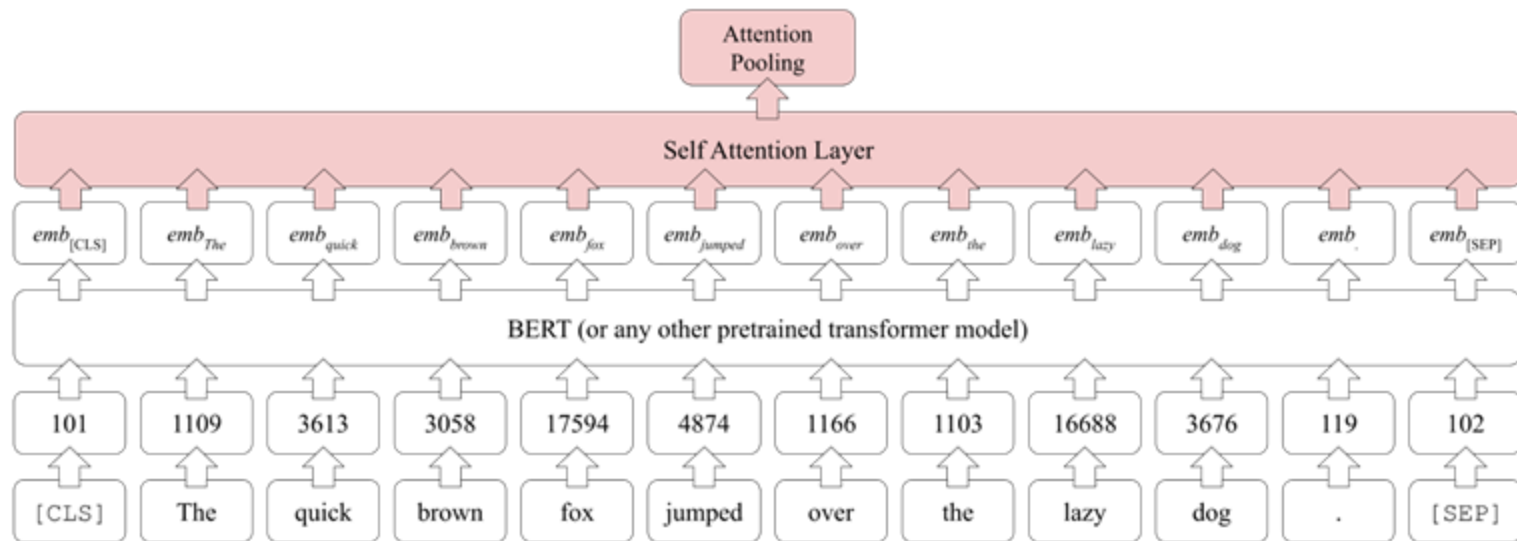
## *Mean Pooling*





# Pooling Methods for Sentence Embeddings

## Attention Pooling



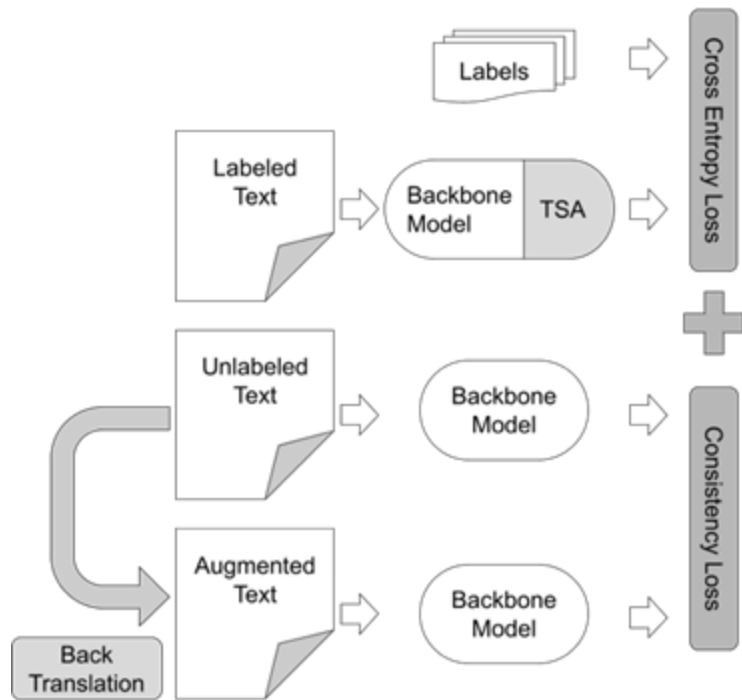
# Semi-Supervised Learning

## *Data Augmentation Techniques*

Original	Once upon a midnight dreary, while I pondered, weak and weary,
Synonym Replacement (EDA)	<b>Erstwhile</b> upon a midnight dreary, while I pondered, weak and weary,
Random Insertion (EDA)	Once upon a midnight dreary, while I pondered, weak and <b>once</b> weary,
Random Swap (EDA)	Once upon <b>I</b> midnight dreary, while <b>a</b> pondered, weak and weary,
Random Delete (EDA)	Once upon a _ dreary, while I pondered, _ and weary,
Back Translation	Once at midnight it was bleak while I was thinking, weak and tired,

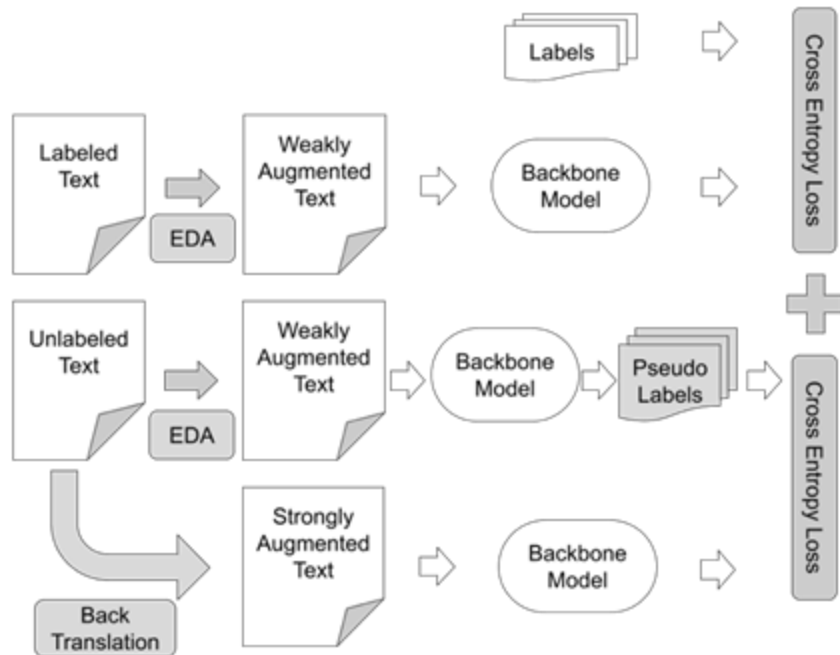
# Semi-Supervised Learning

## *Unsupervised Data Augmentation*



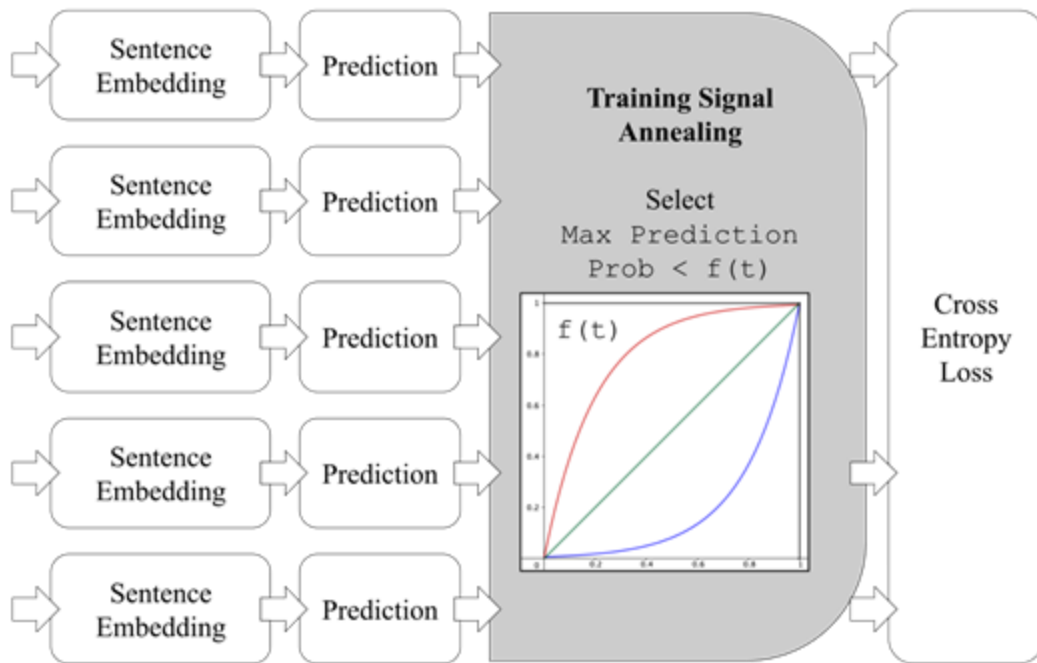
# Semi-Supervised Learning

## *FixMatch*



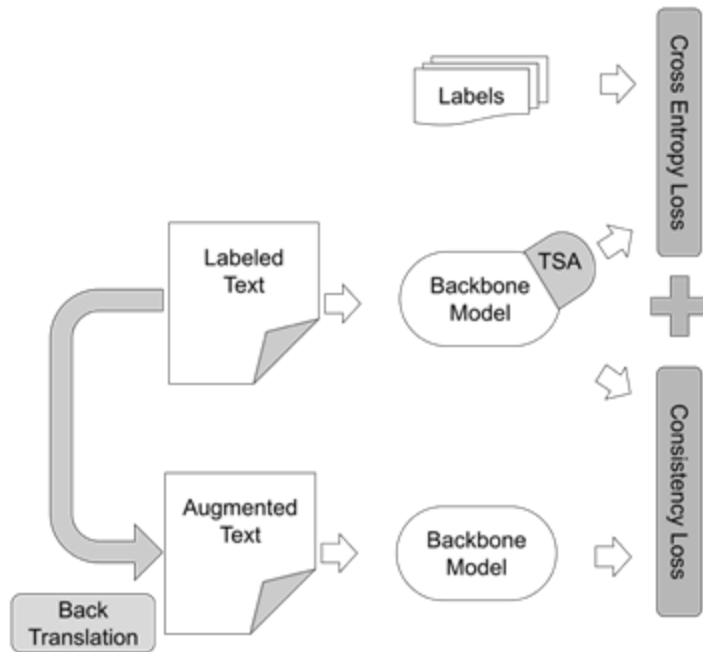
# Loss Engineering

## *Training Signal Annealing*



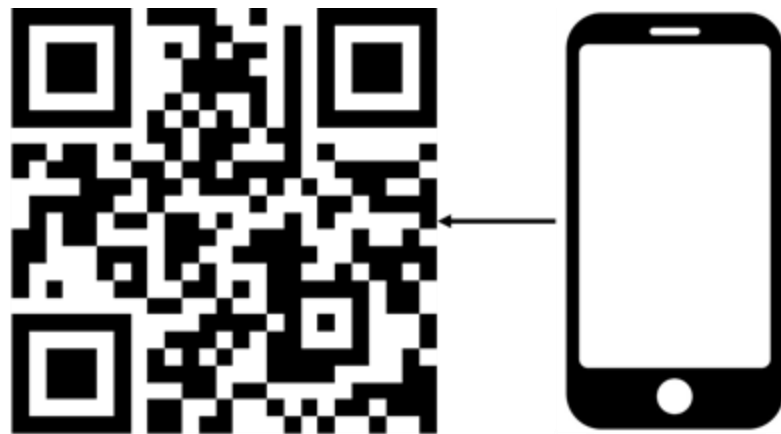
# Loss Engineering

## *Supervised Data Augmentation*



# Results

Model	SectLabel		Extended	
	Macro F1	Micro F1	Macro F1	Micro F1
<i>SciWING</i> (Ramesh Kashyap and Kan, 2020)	0.732	0.900	-	-
RoBERTa-Attn Model (OURS)	0.806	0.904	0.596	0.870
RoBERTa-Attn Model + UDA <sub>log</sub> <sup>†</sup>	0.784	0.906	<b>0.669</b>	<b>0.887</b>
RoBERTa-Attn Model + SDA <sub>log</sub> <sup>†</sup>	<b>0.832</b>	<b>0.929</b>	0.623	0.886
<i>SectLabel</i> (Luong et al., 2010) <sup>‡</sup>	<b>0.847</b>	<b>0.934</b>	-	-



Scan to read the full paper!

Connect with the first author!



[huangpowei@comp.nus.edu.sg](mailto:huangpowei@comp.nus.edu.sg)